

ISSN: 1337-6365

© Slovak University of Technology in Bratislava

All rights reserved

**APLIMAT – JOURNAL OF APPLIED MATHEMATICS**

**VOLUME 2 (2009), NUMBER 2**



# **APLIMAT – JOURNAL OF APPLIED MATHEMATICS**

## **VOLUME 2 (2009), NUMBER 2**

**Edited by:** Slovak University of Technology in Bratislava

**Editor - in - Chief:** KOVÁČOVÁ Monika (Slovak Republic)

**Editorial Board:** CARKOVŠ Jevgenijs (Latvia )  
CZANNER Gabriela (USA)  
CZANNER Silvester (Great Britain)  
DE LA VILLA Augustin (Spain)  
DOLEŽALOVÁ Jarmila (Czech Republic)  
FEČKAN Michal (Slovak Republic)  
FERREIRA M. A. Martins (Portugal)  
FRANCAVIGLIA Mauro (Italy)  
KARPÍŠEK Zdeněk (Czech Republic)  
KOROTOV Sergey (Finland)  
LORENZI Marcella Giulia (Italy)  
MESIAR Radko (Slovak Republic)  
TALAŠOVÁ Jana (Czech Republic)  
VELICHOVÁ Daniela (Slovak Republic)

**Editorial Office:** Institute of natural sciences, humanities and social sciences  
Faculty of Mechanical Engineering  
Slovak University of Technology in Bratislava  
Námestie slobody 17  
812 31 Bratislava

**Correspondence concerning subscriptions, claims and distribution:**

F.X. spol s.r.o  
Azalková 21  
821 00 Bratislava  
journal@aplimat.com

**Frequency:** One volume per year consisting of two issues at price of 120 EUR, per volume,  
including surface mail shipment abroad.  
Registration number EV 2540/08

**Information and instructions for authors are available on the address:** [www.aplimat.com](http://www.aplimat.com)

**Printed by:** F.X. spol s.r.o, Azalková 21, 821 00 Bratislava

**Copyright © STU 2007-2009, Bratislava**

All rights reserved. No part may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without prior written permission from the Editorial Board. All contributions published in the Journal were reviewed with open and blind review forms with respect to their scientific contents.

# APLIMAT – JOURNAL OF APPLIED MATHEMATICS

## VOLUME 2 (2009), NUMBER 2

### DIFFERENTIAL EQUATIONS AND THEIR APPLICATIONS

<b>BAŠTINEC Jaromír, DIBLÍK Josef:</b> OSCILLATION OF SOLUTIONS OF A LINEAR SECOND ORDER DISCRETE DELAYED EQUATION	13
<b>CARKOVA Viktorija, GOLDŠTEINE Jolanta, SWERDAN Myhailo:</b> CONVERGENCE OF LINEAR MARKOV ITERATIONS	19
<b>CARKOVŠ Jevgenijs, EGLE Aigars:</b> ON CONTINUOUS STOCHASTIC MODELING OF HETEROSKEDASTIC CONDITIONAL VARIANCE	29
<b>FAJMON Břetislav, ŠMARDA Zdeněk:</b> APPLICATION OF INTEGRAL INEQUALITIES IN THE THEORY OF INTEGRAL AND INTEGRODIFFERENTIAL EQUATIONS	43
<b>FARAGÓ István, HORVÁTH Róbert, KOTOROV Sergey:</b> DISCRETE MAXIMUM PRINCIPLES FOR PARABOLIC PROBLEMS WITH GENERAL BOUNDARY CONDITIONS	49
<b>FEČKAN Michal:</b> CHAOTIC OSCILLATIONS OF ELASTIC BEAMS	57
<b>FILIPPOVA Olga, ŠMARDA Zdeněk:</b> SINGULAR INITIAL PROBLEM FOR FREDHOLM-VOLTERRA INTEGRODIFFERENTIAL EQUATIONS	69
<b>HERRMANN Leopold, MLS Jiří, ONDOVČIN Tomáš:</b> OSCILLATIONS FOR A HYPERBOLIC DIFFUSION EQUATION WITH TIME-DEPENDENT COEFFICIENTS	77
<b>ILAVSKÁ Iveta, NAJMANOVÁ Anna, OLACH Rudolf:</b> EXISTENCE OF NONOSCILLATORY SOLUTIONS OF NONLINEAR DELAY DIFFERENTIAL EQUATIONS	83
<b>JÁNSKÝ Jiří:</b> THE $\theta$ -METHODS FOR THE DELAY DIFFERENTIAL EQUATIONS	89
<b>KEČKEMÉTYOVÁ Mária, BOCK Igor:</b> A PSEUDOHYPERBOLIC PROBLEM FOR VON A KÁRMÁN SYSTEM	101
<b>MURESAN Anton S.:</b> STRICT FIXED POINT PRINCIPLES AND APPLICATIONS TO MATHEMATICAL ECONOMICS	111
<b>MURESAN Viorica:</b> A FREDHOLM INTEGRAL EQUATION WITH LINEAR MODIFICATION OF THE ARGUMENT	125

<b>PANCZA David:</b> ON ASYMPTOTIC BEHAVIOUR OF THE NONLINEAR VISCOELASTIC MINDLIN-TIMOSHENKO THIN PLATE MODEL	<b>133</b>
<b>ŠMARDA Zdeněk:</b> GENERALIZATION OF CERTAIN INTEGRAL INEQUALITIES	<b>143</b>

# APLIMAT – JOURNAL OF APPLIED MATHEMATICS

## VOLUME 2 (2009), NUMBER 2

### MODELING AND SIMULATION

<b>BARAKEH Bilal:</b> CONVERGENCE PROOF OF A MONTE CARLO SCHEME FOR THE RESOLUTION OF THE SMOLUCHOWSKI COAGULATION EQUATION	149
<b>DAVEAU Christian, KHELIFI Abdessatar, SUSCHENKO Anton:</b> RECONSTRUCTION OF CLOSELY SPACED SMALL INHOMOGENEITIES VIA BOUNDARY MEASUREMENTS FOR THE FULL TIME-DEPENDENT MAXWELL'S EQUATIONS	159
<b>CHATTERJEE Samrat, VENTURINO Ezio:</b> SHARK-FISH INTERPLAY AT DIFFERENT LIFESTAGES	177
<b>JURÍK Tomáš:</b> A NEW ACTIVE-SET METHOD FOR LINEAR PROGRAMMING BASED ON TRANSFORMATION OF FEASIBLE DIRECTION ALGORITHM INTO UNCONSTRAINED MINIMIZATION PROBLEM	189
<b>KINDLER Eugene, KŘIVÝ Ivan:</b> OBJECT-ORIENTED PROGRAMMING LANGUAGES AS TOOLS FOR FORMULATIONS OF SYSTEM ABSTRACTION	197
<b>KISELA Tomáš:</b> FRACTIONAL GENERALIZATION OF THE CLASSICAL VISCOELASTICITY MODELS	209
<b>KVAPIL David:</b> STABILISATION OF MEAN AND VARIANCE FOR NONSTATIONARY PROCESSES	219
<b>LUCKA Maria, PIECKA Stanislav:</b> PARALLEL POSIX THREADS BASED ANT COLONY OPTIMIZATION USING ASYNCHRONOUS COMMUNICATIONS	229
<b>MARČOKOVÁ Mariana, GULDAN Vladimír:</b> ON ONE ORTHOGONAL TRANSFORM APPLIED ON A SYSTEM OF ORTHOGONAL POLYNOMIALS IN TWO VARIABLES	239
<b>PIEKARZ Monika:</b> AN INTROSPECT OF SIMULATION OF NONDETERMINISTIC TURING MACHINE WITH A REAL-ANALYTIC FUNCTION	247





# APLIMAT – JOURNAL OF APPLIED MATHEMATICS

## LIST OF REVIEWERS

<b>Abderraman Jesus C., PhD.</b>	UPM - Technological University of Madrid, Madrid, Spain
<b>Andrade Marina, Dr., PhD.</b>	ISCTE Business School, Lisbon, Portugal
<b>Baranová Eva, RNDr.</b>	Technical university Košice, Košice, Slovak Republic
<b>Beránek Jaroslav, doc. RNDr. CSc.</b>	Masaryk University, Brno, Czech Republic
<b>Daveau Christian, Dr. of Mathematics</b>	Universit'e de Cergy-Pontoise, Cergy-Pontoise Cedex, France
<b>Diblík Josef, Prof. RNDr., DrSc.</b>	Brno University of Technology, Czech Republic
<b>Dobrucky Branislav, Prof.</b>	University of Zilina, Zilina, Slovak Republic
<b>Dorociaková Božena, RNDr., PhD.</b>	University of Zilina, Zilina, Slovak Republic
<b>Došlá Zuzana, Prof. RNDr. DrSc.</b>	Masarykova univerzita, Brno, Czech Republic
<b>Fajmon Břetislav, RNDr., PhD.</b>	Brno University of Technology, Brno, Czech Republic
<b>Ferreira Manuel Alberto M., Full Prof.</b>	ISCTE, Lisboa, Portugal
<b>Filipe Jose Antonio, Assistant Prof.</b>	ISCTE, Lisboa, Portugal
<b>Grodzki Zdzisław, Dr. hab.</b>	Technical University of Lublin, Lublin, Poland
<b>Habiballa Hashim, RNDr. PaedDr., PhD.</b>	University of Ostrava, Ostrava, Czech Republic
<b>Heckenbergerova Jana, Mgr.</b>	University of Pardubice, Pardubice, Czech Republic
<b>Hinterleitner Irena, Mgr.</b>	Brno University of Technology, Brno, Czech Republic
<b>Hošková Šárka, RNDr., PhD.</b>	University of Defence, Brno, Czech Republic
<b>Hřebíček Jiří, Prof. RNDr., CSc.</b>	Masaryk University, Brno, Czech Republic
<b>Hrnčiarová Ľubica, doc.Ing., PhD.</b>	University of Economics in Bratislava, Bratislava, Slovak Republic

<b>Huňka František</b> , doc. Ing., CSc.	Ostravská univerzita in Ostrava, Czech Republic
<b>Chocholatá Michaela</b> , Ing., PhD.	University of Economics, Bratislava, Slovak Republic
<b>Chvalina Jan</b> , Prof. RNDr., DrSc.	Brno University of Technology, Brno, Czech Republic
<b>Chvátalová Zuzana</b> , RNDr., PhD.	Brno University of Technology, Brno, Czech Republic
<b>Jancarik Antonin</b> , PhD.	Charles University, Prague 1, Czech Republic
<b>Kalina Martin</b> , doc. RNDr., CSc.	Slovak University of Technology, Bratislava, Slovak Republic
<b>Kapička Vratislav</b> , Prof. RNDr., DrSc.	Masarykova univerzita, Brno, Czech Republic
<b>Khelifi Abdessatar</b> , PhD.	Sciences of Bizerte, Zarzouna, Tunisia
<b>Kováčová Monika</b> , Mgr., PhD.	Slovak University of Technology, Bratislava, Slovak Republic
<b>Kriz Jan</b> , RNDr., PhD.	University of Hradec Kralove, Hradec Kralove, Czech Republic
<b>Kureš Miroslav</b> , Assoc. Prof.	Brno University of Technology, Brno, Czech Republic
<b>Kvasz Ladislav</b> , doc. Dr., PhD.	Charles University, Prague, Czech Republic
<b>Lovisek Jan</b> , Prof.	Slovak University of Technology, Bratislava, Slovak Republic
<b>Lungu Nicolaie</b> , Prof.	Technical University of Cluj-Napoca, Romania
<b>Malacká Zuzana</b> , RNDr., PhD.	University of Žilina, Slovak Republic
<b>Mamrilla Dušan</b> , doc. RNDr., CSc.	University of Prešov in Prešov, Slovak Republic
<b>Marček Dušan</b>	University of Žilina, Žilina, Slovak Republic
<b>Marčoková Mariana</b> , doc. RNDr., CSc.	University of Žilina, Žilina, Slovak Republic
<b>Maroš Bohumil</b> , doc. RNDr., CSc.	Brno University of Technology, Brno, Czech Republic
<b>Martinek Pavel</b> , Ing., PhD.	Palacky University, Olomouc, Czech Republic
<b>Mikeš Josef</b> , Prof. RNDr., DrSc.	Palacký University, Olomouc, Czech Republic
<b>Mišútová Mária</b> , doc. RNDr., PhD.	Slovak University of Technology, Trnava, Slovak Republic
<b>Moučka Jiří</b> , doc. PhDr., PhD.	University of Defence, Brno, Czech Republic

<b>Pavlačka Ondřej</b> , Mgr., PhD.	Palacký University, Olomouc, Czech Republic
<b>Plháková Alena</b> , doc., PhD.	Palacký University, Olomouc, Czech Republic
<b>Pokorný Milan</b> , PaedDr. PhD.	Trnava University, Trnava, Slovak Republic
<b>Pokorny Michal</b> , Prof.	University of Zilina, Zilina, Slovak Republic
<b>Pribullová Anna</b>	Gephysical Insitute SAS, Bratislava, Slovak Republic
<b>Půlpán Zdeněk</b> , Prof. RNDr., PhD.	University Hradec Králové, Hradec Králové, Czech Republic
<b>Rus Ioan A.</b> , Prof.	Babes-Bolyai University of Cluj-Napoca, Cluj-Napoca, Romania
<b>Slapal Josef</b> , Prof. RNDr., CSc.	Brno University of Technology, Brno, Czech Republic
<b>Svoboda Zdeněk</b> , RNDr., CSc.	Brno University of Technology, Brno, Czech Republic
<b>Šmarda Zdeněk</b> , doc. RNDr., CSc.	Brno University of Technology, Brno, Czech Republic
<b>Talašová Jana</b> , doc. RNDr., CSc.	Palacký University Olomouc, Olomouc, Czech Republic
<b>Vacek Vladimír</b> , RNDr.	Technická univerzita, Zvolen, Slovak Republic
<b>Vajsablova Margita</b> , PhD.	Slovak University of Technology, Bratislava, Slovak Republic
<b>Vanžurová Alena</b> , doc. RNDr., CSc.	Palacký University., Olomouc, Czech Republic
<b>Vávra František</b> , doc. Ing., CSc.	University of West Bohemia, Plzeň, Czech Republic
<b>Velichová Daniela</b> , doc. RNDr., CSc.	Slovak University of Technology, Bratislava, Slovak Republic
<b>Volna Eva</b> , doc. RNDr. PaedDr., PhD.	University of Ostrava, Ostrava 1, Czech Republic
<b>Witkovsky Viktor</b> , doc. RNDr., CSc.	Slovak Academy of Sciences, Bratislava, Slovak Republic
<b>Žáčková Petra</b> , Mgr.	Technical university of Liberec, Liberec, Czech Republic
<b>Žák Libor</b>	University of Technology, Brno, Czech Republic



# OSCILLATION OF SOLUTIONS OF A LINEAR SECOND ORDER DISCRETE DELAYED EQUATION

BAŠTINEC Jaromír, (CZ), DIBLÍK Josef, (CZ)

**Abstract.** A linear second order discrete delayed equation  $\Delta x(n) = -p(n)x(n-1)$  with a positive coefficient  $p$  is considered for  $n \rightarrow \infty$ . This equation is known to have a positive solution if  $p$  fulfils an inequality. The goal of the paper is to show that in the case of an opposite inequality for  $p$  all solutions of the equation considered are oscillating for  $n \rightarrow \infty$ .

**Key words and phrases.** Discrete equation, delay, linear equation, positive solution, oscillating solution.

*Mathematics Subject Classification.* Primary 39A10; Secondary 39A11.

## 1 Introduction

In this remark we consider the delayed scalar linear discrete equation of the second order

$$\Delta x(n) = -p(n)x(n-1) \quad (1)$$

where  $n \in \mathbb{Z}_a^\infty := \{a, a+1, \dots\}$ ,  $a \in \mathbb{N}$  is fixed,  $\Delta x(n) = x(n+1) - x(n)$ ,  $p: \mathbb{Z}_a^\infty \rightarrow \mathbb{R}^+ := (0, \infty)$ .

A solution  $x = x(n): \mathbb{Z}_a^\infty \rightarrow \mathbb{R}$  of (1) is positive on  $\mathbb{Z}_a^\infty$  if  $x(n) > 0$  for every  $n \in \mathbb{Z}_a^\infty$ . A solution  $x = x(n): \mathbb{Z}_a^\infty \rightarrow \mathbb{R}$  of (1) is oscillating on  $\mathbb{Z}_a^\infty$  if it is not positive on  $\mathbb{Z}_{a_1}^\infty$  for arbitrary  $a_1 \in \mathbb{Z}_a^\infty$ .

In the paper [1] a delayed linear difference equation of higher order is considered and the following result related to equation (1) on existence of a positive solution is proved.

**Theorem 1.1** *Let  $a \in \mathbb{N}$ . Suppose that there exists a constant  $\theta \in [0, 1)$  such that the function  $p: \mathbb{Z}_a^\infty \rightarrow \mathbb{R}^+$  satisfies*

$$p(n) \leq \frac{1}{4} + \frac{1}{16n^2} + \frac{\theta}{16(n \ln n)^2} \quad (2)$$

for every  $n \in \mathbb{Z}_a^\infty$ . Then there exists a positive integer  $a_1 \geq a$  and a solution  $u = u(k)$ ,  $k \in N(a_1)$  of equation (1) such that the inequalities

$$0 < x(n) < \left(\frac{1}{2}\right)^n \cdot \sqrt{n \ln n}$$

hold for every  $n \in \mathbb{Z}_{a_1}^\infty$ .

Our goal is to answer an open question formulated in [1], whether all solutions of (1) are oscillating if inequality (2) is replaced by an opposite inequality

$$p(n) \geq \frac{1}{4} + \frac{1}{16n^2} + \frac{\kappa}{16(n \ln n)^2} \quad (3)$$

assuming  $\kappa \geq 1$  and  $n$  is sufficiently large.

From recent investigation performed in [2] follows that Theorem 1.1 holds even if  $\theta = 1$  and consequently the answer is negative if we admit  $\kappa = 1$  and equality instead of inequality in (3).

Below we prove that if the inequality in (3) is strong and  $\kappa > 1$ , then the answer is positive - all solutions are oscillatory.

The proof of our main result will use a consequence from one of results of Y. Domshlak [6, p. 69].

**Lemma 1.2** *Let  $q$  and  $r$  be arbitrary natural numbers such that  $r - q > 1$ . Let  $\{\varphi(n)\}_1^\infty$  be a given sequence of positive numbers and  $\nu_0$  be a positive number such that there exists a positive number  $\nu = \nu(q, r) < \nu_0$  satisfying*

$$\sum_{q+1}^r \varphi(n) \leq \frac{\pi}{\nu}, \quad \frac{\pi}{\nu} \leq \sum_{q+1}^{r+1} \varphi(n) \leq \frac{2\pi}{\nu}. \quad (4)$$

Then, if  $p(q+1) \geq 0$  and for  $n \in \mathbb{Z}_{q+2}^r$

$$p(n) \geq \frac{\sin \nu \varphi(n-1) \cdot \sin \nu \varphi(n+1)}{\sin \nu [\varphi(n-1) + \varphi(n)] \cdot \sin \nu [\varphi(n) + \varphi(n+1)]} \quad (5)$$

holds, then all solutions of the equation

$$x(n+1) - x(n) + p(n)x(n-1) = 0$$

are oscillatory.

Moreover, we will use an auxiliary result giving the asymptotic decomposition of the logarithm (cf. e.g. [5]). The symbols “ $o$ ” and “ $O$ ” used below stands for the Landau order symbols.

**Lemma 1.3** *For fixed  $r, \sigma \in \mathbb{R} \setminus \{0\}$  the asymptotic representation*

$$\frac{\ln^\sigma(n-r)}{\ln^\sigma n} = 1 - \frac{r\sigma}{n \ln n} - \frac{r^2\sigma}{2n^2 \ln n} + \frac{r^2\sigma(\sigma-1)}{2(n \ln n)^2} + o\left(\frac{1}{n^3}\right)$$

holds for  $n \rightarrow \infty$ .

## 2 Main Result

In this part we give sufficient conditions for all solutions of equation (1) to be oscillatory when  $n \rightarrow \infty$ .

**Theorem 2.1** *Let  $a \in \mathbb{N}$  be sufficiently large and  $\kappa > 1$ . Suppose that the function  $p: \mathbb{Z}_a^\infty \rightarrow \mathbb{R}^+$  satisfies*

$$p(n) > \frac{1}{4} + \frac{1}{16n^2} + \frac{\kappa}{16(n \ln n)^2} \quad (6)$$

*for every  $n \in \mathbb{Z}_a^\infty$ . Then all solutions of equation (1) are oscillating when  $n \rightarrow \infty$ .*

**Proof.** We set

$$\varphi(n) = \frac{1}{n \ln n}$$

and consider the asymptotic decomposition of  $\varphi(n-1)$  when  $n$  is sufficiently large. Applying Lemma 1.3 (for  $\sigma = -1$  and  $\tau = 1$ ) we get

$$\begin{aligned} \varphi(n-1) &= \frac{1}{(n-1) \ln(n-1)} = \frac{1}{n(1-1/n) \ln(n-1)} \\ &= \frac{1}{n \ln n} \left( 1 + \frac{1}{n} + \frac{1}{n^2} + O\left(\frac{1}{n^3}\right) \right) \left( 1 + \frac{1}{n \ln n} + \frac{1}{2n^2 \ln n} - \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right) \right) \\ &= \frac{1}{n \ln n} \left( 1 + \frac{1}{n} + \frac{1}{n \ln n} + \frac{1}{n^2} + \frac{3}{2n^2 \ln n} + \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right) \right). \end{aligned}$$

Similarly, applying Lemma 1.3 (for  $\sigma = -1$  and  $\tau = -1$ ) we obtain

$$\begin{aligned} \varphi(n+1) &= \frac{1}{(n+1) \ln(n+1)} = \frac{1}{n(1+1/n) \ln(n+1)} \\ &= \frac{1}{n \ln n} \left( 1 - \frac{1}{n} + \frac{1}{n^2} + O\left(\frac{1}{n^3}\right) \right) \left( 1 - \frac{1}{n \ln n} + \frac{1}{2n^2 \ln n} + \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right) \right) \\ &= \frac{1}{n \ln n} \left( 1 - \frac{1}{n} - \frac{1}{n \ln n} + \frac{1}{n^2} + \frac{3}{2n^2 \ln n} + \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right) \right). \end{aligned}$$

Next, we develop asymptotic decomposition for  $\varphi(n-1)\varphi(n+1)$  when  $n$  is sufficiently large. Using previous decompositions we have

$$\begin{aligned} \varphi(n-1)\varphi(n+1) &= \frac{1}{(n \ln n)^2} \\ &\quad \times \left( 1 + \frac{1}{n} + \frac{1}{n \ln n} + \frac{1}{n^2} + \frac{3}{2n^2 \ln n} + \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right) \right) \\ &\quad \times \left( 1 - \frac{1}{n} - \frac{1}{n \ln n} + \frac{1}{n^2} + \frac{3}{2n^2 \ln n} + \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right) \right) \\ &= \frac{1}{(n \ln n)^2} \left( 1 + \frac{1}{n^2} + \frac{1}{n^2 \ln n} + \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right) \right). \end{aligned}$$

Recalling the asymptotical decomposition of  $\sin x$  when  $x \rightarrow 0$ :  $\sin x = x + O(x^3)$  we get (preserving the order of accuracy  $O(n^{-3})$ )

$$\begin{aligned}\sin \nu \varphi(n-1) &= \nu \varphi(n-1) + O\left(\frac{\nu^3}{n^3}\right), \\ \sin \nu \varphi(n+1) &= \nu \varphi(n+1) + O\left(\frac{\nu^3}{n^3}\right), \\ \sin \nu [\varphi(n-1) + \varphi(n)] &= \nu [\varphi(n-1) + \varphi(n)] + O\left(\frac{\nu^3}{n^3}\right), \\ \sin \nu [\varphi(n) + \varphi(n+1)] &= \nu [\varphi(n) + \varphi(n+1)] + O\left(\frac{\nu^3}{n^3}\right)\end{aligned}$$

when  $n \rightarrow \infty$ . Then it is easy to see that for the right-hand side of the inequality (5) we have

$$\begin{aligned}\mathcal{R} &:= \frac{\sin \nu \varphi(n-1) \cdot \sin \nu \varphi(n+1)}{\sin \nu [\varphi(n-1) + \varphi(n)] \cdot \sin \nu [\varphi(n) + \varphi(n+1)]} \\ &= \frac{\left(\nu \varphi(n-1) + O\left(\frac{\nu^3}{n^3}\right)\right) \cdot \left(\nu \varphi(n+1) + O\left(\frac{\nu^3}{n^3}\right)\right)}{\left(\nu [\varphi(n-1) + \varphi(n)] + O\left(\frac{\nu^3}{n^3}\right)\right) \cdot \left(\nu [\varphi(n) + \varphi(n+1)] + O\left(\frac{\nu^3}{n^3}\right)\right)} \\ &= \mathcal{R}_1 + O\left(\frac{\nu^2}{n^3}\right), \quad n \rightarrow \infty\end{aligned}$$

where

$$\mathcal{R}_1 := \frac{\varphi(n-1)\varphi(n+1)}{\varphi^2(n) + \varphi(n)\varphi(n-1) + \varphi(n)\varphi(n+1) + \varphi(n-1)\varphi(n+1)}.$$

Moreover, for  $\mathcal{R}_1$  we get an asymptotical decomposition when  $n \rightarrow \infty$ . We represent  $\mathcal{R}_1$  in the form

$$\begin{aligned}\mathcal{R}_1 &= \frac{\varphi(n-1)\varphi(n+1)}{\varphi^2(n) + \varphi(n)\varphi(n-1) + \varphi(n)\varphi(n+1) + \varphi(n-1)\varphi(n+1)} \\ &= \frac{\frac{\varphi(n-1)\varphi(n+1)}{\varphi^2(n)}}{1 + \frac{\varphi(n-1)}{\varphi(n)} + \frac{\varphi(n+1)}{\varphi(n)} + \frac{\varphi(n-1)\varphi(n+1)}{\varphi^2(n)}}.\end{aligned}$$

Because

$$\begin{aligned}\frac{\varphi(n-1)\varphi(n+1)}{\varphi^2(n)} &= 1 + \frac{1}{n^2} + \frac{1}{n^2 \ln n} + \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right), \\ \frac{\varphi(n-1)}{\varphi(n)} &= 1 + \frac{1}{n} + \frac{1}{n \ln n} + \frac{1}{n^2} + \frac{3}{2n^2 \ln n} + \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right), \\ \frac{\varphi(n+1)}{\varphi(n)} &= 1 - \frac{1}{n} - \frac{1}{n \ln n} + \frac{1}{n^2} + \frac{3}{2n^2 \ln n} + \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right),\end{aligned}$$



$$\begin{aligned}
\mathcal{R}_1 &= \left(1 + \frac{1}{n^2} + \frac{1}{n^2 \ln n} + \frac{1}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right)\right) \\
&\quad \times \left(4 + \frac{3}{n^2} + \frac{4}{n^2 \ln n} + \frac{3}{(n \ln n)^2} + O\left(\frac{1}{n^3}\right)\right)^{-1} \\
&= \frac{1}{4} \left(1 + \frac{1}{n^2} + \frac{1}{n^2 \ln n} + \frac{1}{(n \ln n)^2}\right) \left(1 - \frac{3}{4n^2} - \frac{1}{n^2 \ln n} - \frac{3}{4(n \ln n)^2}\right) + O\left(\frac{1}{n^3}\right) \\
&= \frac{1}{4} \left(1 + \frac{1}{4n^2} + \frac{1}{4(n \ln n)^2}\right) + O\left(\frac{1}{n^3}\right).
\end{aligned}$$

Finalizing our decompositions we see that

$$\mathcal{R} = \frac{1}{4} \left(1 + \frac{1}{4n^2} + \frac{1}{4(n \ln n)^2}\right) + O\left(\frac{1}{n^3}\right) + O\left(\frac{\nu^2}{n^3}\right).$$

It is easy to see, that the inequality (5) turns into

$$p(n) \geq \frac{1}{4} \left(1 + \frac{1}{4n^2} + \frac{1}{4(n \ln n)^2}\right) + O\left(\frac{1}{n^3}\right) + O\left(\frac{\nu^2}{n^3}\right). \quad (7)$$

We conclude that if  $p(n)$  satisfies (6) and  $\nu \in (0, \nu_0]$  with  $\nu_0$  fixed, then there exists an index  $n_{\nu_0}$  such that for  $n \geq n_{\nu_0}$  inequality (7) is satisfied. Then the assumption (5) of Lemma 1.2 holds. Inequalities (4) are valid as well because  $\nu$  can be taken sufficiently small. This fact ends the proof.

**Remark 2.2** For further reading we recommend relevant literature, e.g. book [8] and, except the above mentioned, papers [3, 4], [7], [9]–[14].

## Acknowledgement

The first author was supported by grant 201/08/0469 of the Czech Grant Agency (Prague) and by the Council of Czech Government MSM 0021630529. The second author was supported by grant 201/07/0145 of the Czech Grant Agency (Prague) and by the Council of Czech Government MSM 00216 30503 and MSM 00216 30519.

## References

- [1] BAŠTINEC, J., DIBLÍK, J.: *Remark on positive solutions of discrete equation  $\Delta u(k+n) = -p(k)u(k)$* , Nonlinear Anal., **63** (2005), e2145–e2151.
- [2] BAŠTINEC, J., DIBLÍK, J., ŠMARDA, Z.: *Existence of Positive Solutions of Discrete Linear Equations with a Single Delay*, J. Difference Equ. Appl., submitted.
- [3] BEREZANSKY, L., BRAVERMAN, E.: *On existence of positive solutions for linear difference equations with several delays*, Adv. Dyn. Syst. Appl., **1** (2006), 29–47.

- [4] DIBLÍK, J. *Asymptotic behaviour of solutions of discrete equations*, Funct. Differ. Equ., **11** (2004), 37–48.
- [5] DIBLÍK, J., KOKSCH, N.: *Positive solutions of the equation  $\dot{x}(t) = -c(t)x(t - \tau)$  in the critical case*, J. Math. Anal. Appl., **250** (2000), 635–659.
- [6] DOMSHLAK, Y.: *Oscillation properties of discrete difference inequalities and equations: The new approach*, Funct. Differ. Equ., **1**, 60–82 (1993).
- [7] DOMSHLAK, Y., STAVROULAKIS, I.P.: *Oscillation of first-order delay differential equations in a critical case*, Appl. Anal., **61** (1996), 359–371.
- [8] GYÖRI, I., LADAS, G.: *Oscillation Theory of Delay Differential Equations*, Clarendon Press (1991).
- [9] GYÖRI, I., PITUK, M.: *Asymptotic formulae for the solutions of a linear delay difference equation*, J. Math. Anal. Appl., **195** (1995), 376–392.
- [10] KARAJANI, P., STAVROULAKIS, I.P.: *Oscillation criteria for second-order delay and functional equations*, Stud. Univ. Žilina, Math. Ser., **18**, No 1, (2004), 17–26.
- [11] KIKINA, L.K., STAVROULAKIS, I.P.: *A survey on the oscillation of solutions of first order delay difference equations*, Cubo, **7**, No 2, (2005), 223–236.
- [12] MIGDA, M., MIGDA, J.: *Asymptotic behaviour of solutions of difference equations of second order*, Demonstr. Math., **XXXII** (1999), 767–773.
- [13] STAVROULAKIS, I.P.: *Oscillation criteria for first order delay difference equations*, Mediterr. J. Math., **1**, No 2, (2004), 231–240.
- [14] STAVROULAKIS, I.P.: *Oscillation criteria for delay and difference equations*, Stud. Univ. Žilina, Math. Ser., **17**, No 1, (2003), 161–167.

### **Current address**

#### **Baštinec Jaromír, doc. RNDr., CSc.**

Department of Mathematics  
Faculty of Electrical Engineering and Communication  
Brno University of Technology  
Technická 8  
616 00 Brno  
Czech Republic  
tel.: +420-541143222  
email: bastinec@feec.vutbr.cz

#### **Diblík Josef, Prof. RNDr., DrSc.**

Brno University of Technology  
Czech Republic  
tel.: +420-514147601, +420-541143155  
email: diblik.j@fce.vutbr.cz, diblik@feec.vutbr.cz

## CONVERGENCE OF LINEAR MARKOV ITERATIONS

CARKOVA Viktorija, (LV), GOLDŠTEINE Jolanta, (LV),  
SWERDAN Myhailo, (UA)

**Abstract.** This paper deals with the linear difference equation in  $\mathbb{R}^n$  with coefficients dependent on the Markov chain. It is proved that covariance matrices of solutions can be analyzed using powers of a positive operator in a Banach space with a reproducing cone. This property permits to formulate the necessary and sufficient mean square stability condition as a spectral problem or a problem of positive solvability of a specially constructed linear operator equation. The paper discusses three possible approaches for convergence analysis of defined by difference equation iterative procedure: mean square stability analysis using the second Lyapunov method; mean square stability analysis using Lyapunov index; reducibility method for moments of solutions, which permits approximately to reduce mean square stability problem to analysis of equation with constant coefficients.

**Key words and phrases.** Stochastic difference equations, Markov dynamical systems, mean square stability.

*Mathematics Subject Classification.* Primary 60A05, 08A72; Secondary 28E10.

## 1 Introduction

The paper deals with asymptotic stability problem for linear difference equations with Markov coefficients. A linear difference equation in  $\mathbb{R}^n$  defined by equality:

$$x_t = A(y_t)x_{t-1}, t \in \mathbb{N}, \quad (1)$$

where  $\{A(y), y \in \mathbb{Y}\}$  is continuous  $n \times n$  matrix function on the metric compact  $\mathbb{Y}$ ,  $\sup_y \|A(y)\| = \text{const} < \infty$ ;  $\{y_t, t \in \mathbb{N}\}$  is a homogeneous Feller Markov chain with phase space  $\mathbb{Y}$ , invariant measure  $\mu(dy)$ , and transition probability  $p(y, dz)$ . Under initial conditions  $x_k = x$ ,  $y_k = y$  the

vector function  $x_t(k, x, y) = X(t, k, y)x$ , where  $X(t, k, y) := \prod_{m=k+1}^t A(y_m)$ , satisfies the difference equation (1) for any  $t \geq k$ .

This paper discusses three possible approaches for convergence analysis of iterative procedure (1) with Markov coefficients:

- mean square stability analysis of (1) using the second Lyapunov method;
- mean square stability analysis using Lyapunov index for (1);
- reducibility method for moments of equation (1) solutions, which permits to reduce (1) to equation with constant coefficients.

## 2 Mean square stability analysis using the second Lyapunov method

The most powerful tool for asymptotic stability analysis of dynamical system is the Second Lyapunov method. One should choose a nonnegative function  $V(x, y)$ , satisfying an equality  $V(x, y) = 0$  if and only if  $x = 0$  (called Lyapunov function) and analyze an expectation of difference by virtue of the above system and Markov chain (or the Lyapunov operator  $\mathbf{L}$ ) defined by equality  $(\mathbf{L}V)(x, y) := \mathbb{E}\{V(x_t, y_t)/x_{t-1} = x, y_{t-1} = y\} - V(x, y)$ . If there exists such Lyapunov function that  $|x|^p \leq V(x, y) \leq c_1|x|^p$ ,  $(\mathbf{L}V)(x, y) \leq -c_2|x|^p$  for any  $x, y$  and some positive  $p, c_1, c_2$  then with increasing of  $t$  to infinity any solution of the above difference equations exponentially tends to zero with probability one. The main idea of this approach is to choose for (1) the Lyapunov function as a quadratic form  $V(x, y) := (v(y)x, x)$  and then to analyze spectral properties of linear operator defined by an equality  $((\mathbf{A}v)(y)x, x) := (\mathbf{L}V)(x, y)$ .

**Definition 2.1** *The equation (1) is called as exponentially mean square stable if there exist such constants  $c > 0$  and  $\lambda \in (0, 1)$  that*

$$\mathbb{E}|x_t(k, x, y)|^2 \leq c\lambda^{t-k}|x|^2 \quad (2)$$

for any  $y \in \mathbb{Y}$ ,  $x \in \mathbb{R}^n$ ,  $k \in \mathbb{N}$  and  $t \geq k$ .

Let  $\mathbb{V}$  be the Banach space of symmetric uniformly bounded continuous  $n \times n$  matrix functions  $\{q(y), y \in \mathbb{Y}\}$  with norm

$$\|q\| := \sup_{y \in \mathbb{Y}, \|x\|=1} |(q(y)x, x)|.$$

To derive mean square stability conditions for (1) special constructed operator equation for quadratic functionals  $(q(y)x, x)$  with  $q \in \mathbb{V}$  is used, where  $(\cdot, \cdot)$  denotes a scalar product. Using matrix  $A(y)$  and transition probability one can define on  $\mathbb{V}$  the linear continuous operator

$$(\mathbf{A}q)(y) := \int_{\mathbb{Y}} A^T(z)q(z)A(z)p(y, dz), \quad (3)$$

where top index  $T$  denotes transposition. It is easy to see that the above defined operator leaves as invariant the cone [5]

$$\mathbb{K} := \{q \in \mathbb{V} : \inf_{y \in \mathbb{Y}, \|x\|=1} (q(y)x, x) \geq 0\}$$

with a set of inner points

$$\overset{\circ}{\mathbb{K}} := \{q \in \mathbb{V} : \inf_{y \in \mathbb{Y}, \|x\|=1} (q(y)x, x) > 0\}.$$

This cone permits to put space  $\mathbb{V}$  in partial order using "inequality"  $q_1 \ll q_2$  if  $q_2 - q_1 \in \mathbb{K}$ . Obviously that  $q \in \overset{\circ}{\mathbb{K}}$  if and only if there exists a such positive constant  $c(q)$  that  $q \gg c(q)I$  where  $I$  is the matrix unit of the space  $\mathbb{V}$ .

**Lemma 2.2** For any  $q \in \mathbb{V}$ ,  $t > k \geq 0$ ,  $y \in \mathbb{Y}$ , and  $x \in \mathbb{R}^n$

$$((\mathbf{A}^t q)(y)x, x) = \mathbb{E} \{(q(y_{t+k})x_{t+k}(k, x, y), x_{t+k}(k, x, y))/y_k = y\}.$$

Using the definition of Cauchy matrix family one can rewrite the assertion of Lemma 1 in the matrix form

$$(\mathbf{A}^t q)(y) = \mathbb{E} \{X^T(t+k, k, y)q(y_{t+k})X(t+k, k, y)/y_k = y\}. \quad (4)$$

**Theorem 2.3** The next assertions are equivalent:

(i) equation (1) is exponentially mean square stable;

(ii) there exists such  $q \in \overset{\circ}{\mathbb{K}}$  that

$$\mathbf{A}q - q = -I; \quad (5)$$

(iii) maximal positive real spectrum point  $\mathbf{r}\{\mathbf{A}\}$  of operator  $\mathbf{A}$  is less than one.

**Proof.** (i)  $\longrightarrow$  (ii). On the basis of an equality

$$\|\mathbb{E} \{X^T(t, 0, y)X(t, 0, y)\}\| = \sup_{|x|=1} |(\mathbb{E} \{X^T(t, 0, y)X(t, 0, y)\}x, x)| =$$

$$\sup_{|x|=1} |\mathbb{E} \{(X(t, 0, y)x, X(t, 0, y)x)\}| = \sup_{|x|=1} \mathbb{E} \{|x_t(0, x, y)|^2\}$$

and mean square stability of (1) there exists the matrix function defined by formula

$$q(y) := \sum_{t=0}^{\infty} \mathbb{E} \{X^T(t, 0, y)X(t, 0, y)\}$$

Because of identity  $X(k, k, y) \equiv I$  and equality

$$\sum_{t=0}^{\infty} \mathbb{E} \{X^T(t, 0, y)X(t, 0, y)\} = I + \sum_{t=1}^{\infty} \mathbb{E} \{X^T(t, 0, y)X(t, 0, y)\}$$

one can write inequality  $q \gg I$ . Therefore  $q \in \overset{\circ}{\mathbb{K}}$ . To complete a proof of the first assertion one can apply formula (4) with matrix function  $q(y) \equiv I$  and to write the equalities

$$\begin{aligned} \mathbf{A}q(y) - q(y) &= \mathbf{A} \left( \sum_{t=0}^{\infty} \mathbf{A}^t I \right) - \sum_{t=0}^{\infty} \mathbf{A}^t I = \\ &= \sum_{t=0}^{\infty} \mathbf{A}^{t+1} I - \sum_{t=0}^{\infty} \mathbf{A}^t I = -I. \end{aligned}$$

(ii)  $\longrightarrow$  (iii). Let  $q \in \overset{\circ}{\mathbb{K}}$  satisfies (5). There exists such positive number  $c(q)$  that  $c(q)I \ll q \ll \|q\|I$  and one can get from the equation (5) inequality  $\mathbf{A}q - q \ll -q/\|q\|$  or  $\mathbf{A}^t q \ll r^t q$  for any  $t \in \mathbb{N}$  where  $r = 1 - \|q\|^{-1} \in (0, 1)$ . Therefore

$$\mathbf{A}^t I \ll \frac{1}{c(q)} \mathbf{A}^t q \ll \frac{r^t}{c(q)} q \ll \|q\| \frac{r^t}{c(q)} I$$

for any  $t \in \mathbb{N}$ , t.i.

$$\sum_{t=0}^m \mathbf{A}^t I \ll \frac{\|q\|}{c(q)} \sum_{t=0}^m r^t I \ll \frac{\|q\|}{c(q)(1-r)} I$$

for any  $m \in \mathbb{N}$  and

$$\lim_{m \rightarrow \infty} \sup_{|x|=1, y \in \mathbb{Y}} \sum_{t=0}^m |((\mathbf{A}^t g)(y)x, x)| < \infty \quad (6)$$

for any  $g \in \mathbb{V}$ . Because linear operator  $\mathbf{A}$  leaves the solid cone  $\mathbb{K}$  as invariant there exists [5] such real spectrum point  $\rho(\mathbf{A})$  that  $\rho(\mathbf{A}) = \sup\{|z|, z \in \sigma(\mathbf{A})\}$  and real eigenfunction  $q_\rho \in \mathbb{K}$  corresponding to this spectrum point, t.i.  $\mathbf{A}q_\rho = \rho(\mathbf{A})q_\rho$ . Therefore if  $\rho(\mathbf{A}) \geq 1$  one should write

$$\lim_{m \rightarrow \infty} \sup_{|x|=1, y \in \mathbb{Y}} \sum_{t=0}^m |((\mathbf{A}^t q_\rho)(y)x, x)| = \infty.$$

This equality contradicts to (6).

(iii)  $\longrightarrow$  (i). Because operator  $\mathbf{A}$  leaves the above defined cone  $\mathbb{K}$  as invariant, there exists [5] positive spectrum point  $\mathbf{r}(\mathbf{A})$  satisfying equality  $\mathbf{r}(\mathbf{A}) = \max \mathbf{Re}\{\sigma(\mathbf{A})\}$ . Therefore, if  $\mathbf{r}(\mathbf{A}) < 1$  then  $\sigma(\mathbf{A}) \subset \{z \in \mathbb{C} : |z| < 1\}$  and there exist [5] such constants  $c > 0$ ,  $\lambda \in (0, 1)$  that  $\|\mathbf{A}^t\| \leq c\lambda^t$  for any  $t \in \mathbb{N}$ . Now one can write inequality

$$\mathbb{E}|x_{t+k}(k, x, y)|^2 = ((\mathbf{A}^t I)(y)x, x) \leq c\lambda^t |x|^2$$

and proof is complete.

More simple stability criterion one can reach assuming that the sequence  $\{y_t, t \in \mathbb{N}\}$  consists of independent random variables with the same distribution  $p(dy)$ . In this case we will consider a contraction  $\hat{\mathbf{A}}$  of the defined by (3) operator  $\mathbf{A}$  on the space  $\mathbb{M}_n$  of symmetric  $n \times n$  real matrices

$$\hat{\mathbf{A}}q := \mathbb{E}\{A^T(y_t)qA(y_t)\} = \int_{\mathbb{Y}} A^T(y)qA(y)p(dy)$$

Using cone of the positive defined matrices  $\overset{\circ}{\mathbb{K}}_n := \mathbb{M}_n \cap \overset{\circ}{\mathbb{K}}$

**Corollary 2.4** *If the sequence  $\{y_t, t \in \mathbb{N}\}$  consists of independent random variables with the same distribution  $p(dy)$  the next assertions are equivalent:*

(i) *equation (1) is exponentially mean square stable;*

(ii) *there exists such  $q \in \overset{\circ}{\mathbb{K}}_n$  that*

$$\hat{\mathbf{A}}q - q = -I;$$

(iii) *maximal positive real spectrum point  $\mathbf{r}(\hat{\mathbf{A}})$  of operator  $\hat{\mathbf{A}}$  is less than one.*

### 3 Mean square Lyapunov index for Markov iterations

If (1) is equation with near to constant coefficients, i.e. matrix in the right part of equation (1) has a form

$$A(y, \varepsilon) = A_0 + \varepsilon A_1(y) + \varepsilon^2 A_2(y) + \dots, \quad (7)$$

the paper proposes an algorithm, which reduces the performances of the equation (1) second moments dynamics to analysis of the operator  $\mathbf{A}(\varepsilon)$  in finite dimensional subspace  $\mathbb{V}(\varepsilon) \subset \mathbb{V}$ . This subspace as well as the restriction matrix  $\Lambda(\varepsilon)$  of the operator  $\mathbf{A}$  on it may be defined by the specially constructed basis  $\mathbf{B}(\varepsilon)$ , analytically dependent on  $\varepsilon$ . The maximal by modulus real eigenvalue  $\rho(\varepsilon)$  of matrix  $\Lambda(\varepsilon)$  for sufficiently small  $\varepsilon > 0$  coincides with similar eigenvalue of operator  $\mathbf{A}(\varepsilon)$ . By terminology of [2] this number defines mean square Lyapunov index by formula  $\lambda_2(\varepsilon) = \lim_{t \rightarrow \infty} \sup_{k, y, |x|=1} \frac{1}{2t} \ln \mathbb{E}\{|x_t(k, x, y)|^2\}$  and this number defines behavior of the second moment  $\mathbb{E}\{|x_t(k, x, y)|^2\}$  as  $t \rightarrow \infty$ : if  $\lambda_2(\varepsilon) < 0$  sequence  $\mathbb{E}\{|x_t(k, x, y)|^2\}$  exponentially decreases, if  $\lambda_2(\varepsilon) > 0$  - exponentially increases.

Let  $\sigma(\mathbf{A})$  be the spectrum and  $r(\mathbf{A})$  be the spectral radius of operator  $\mathbf{A}$ . Substituting matrix (7) in formula (3) one can decompose the operator family  $\mathbf{A}(\varepsilon)$  by power of  $\varepsilon$ :  $\mathbf{A}(\varepsilon) = \mathbf{A}_0 + \varepsilon \mathbf{A}_1 + \varepsilon^2 \mathbf{A}_2 + \dots$  with some bounded operators  $\mathbf{A}_k, k = 1, 2, \dots$  and  $\mathbf{A}_0 q := \int_{\mathbb{Y}} A_0^T q(z) A_0 p(y, dz)$ .

It means that operator family  $\mathbf{A}(\varepsilon)$  analytically depends on parameter  $\varepsilon$  and for finding mean square Lyapunov index  $\lambda_2(\varepsilon)$  we can apply methods and results of perturbation theory of linear continuous operators [4] for decomposition of finite dimension spectral point  $r(\mathbf{A}(\varepsilon))$ . Using the definition of the operator  $\mathbf{A}_0$  we can write that  $\sigma(\mathbf{A}_0) = \{\lambda_1 \cdot \lambda_2 \cdot \lambda_3 : \lambda_{1,2} \in \sigma(A_0), \lambda_3 \in \sigma(P)\}$ . According to this formula spectral radius of operator  $\mathbf{A}_0$  is spectral point which corresponds  $r(\mathbf{A}_0) = \{\max |\lambda|^2 : \lambda \in \sigma(A_0)\}$ , and besides  $r(\mathbf{A}_0) \in \sigma(\mathbf{A}_0)$ . Owing to analyticity of operator-family  $\mathbf{A}(\varepsilon)$  for sufficiently small values of  $\varepsilon$  there exists [4] part of spectrum  $\sigma_\varepsilon$  of the operator  $\mathbf{A}(\varepsilon)$  satisfying the equality  $\lim_{\varepsilon \rightarrow 0} \sigma_\varepsilon = \{r^2(A_0)\}$  where  $r(A_0) = \max\{|\lambda| : \lambda \in \sigma(A_0)\}$ . The root subspace  $\mathbb{V}(\varepsilon) \subset \mathbb{V}$  corresponding to the part of spectrum given by the above formula has the same dimension  $m = \dim \mathbb{V}(0)$  for all sufficiently small  $\varepsilon \geq 0$  [4]. A basis  $\mathbf{B}(\varepsilon)$  can be constructed in  $\mathbb{V}(\varepsilon)$  [4] of the form  $\mathbf{B}(\varepsilon) = \mathbf{P}(\varepsilon) \mathbf{B}^0$ , where  $\mathbf{P}(\varepsilon)$  is the total projector in  $\mathbb{V}(\varepsilon)$  and  $\mathbf{B}^0 \subset \hat{M}(\mathbb{R}^n)$ , where  $\hat{M}(\mathbb{R}^n)$  is a set of symmetric  $n \times n$  matrices, because all corresponding to  $r(\mathbf{A}_0)$  eigen-elements of the operator  $\mathbf{A}_0$  are symmetric  $n \times n$  matrices. Because the projector  $\mathbf{P}(\varepsilon)$  is an analytic function of  $\varepsilon$  [4] one can look for the basis as decomposition

$\mathbf{B}(\varepsilon) = \mathbf{B}^0 + \varepsilon\mathbf{B}^1 + \varepsilon^2\mathbf{B}^2 + \dots$ . Let  $\Lambda(\varepsilon)$  be the matrix of restriction of the operator  $\mathbf{A}(\varepsilon)$  on the subspace  $\mathbb{V}(\varepsilon)$ . Consequently this matrix can be obtained [4] from the expression

$$\mathbf{A}(\varepsilon)\mathbf{B}(\varepsilon) = \mathbf{B}(\varepsilon)\Lambda(\varepsilon) \quad (8)$$

where for the matrix  $\Lambda(\varepsilon)$  also can be used the decomposition  $\Lambda(\varepsilon) = \Lambda_0 + \varepsilon\Lambda_1 + \varepsilon^2\Lambda_2 + \dots$ . Therefore (8) can be rewritten into the form

$$(\mathbf{A}_0 + \varepsilon\mathbf{A}_1 + \varepsilon^2\mathbf{A}_2 + \dots)(\mathbf{B}^0 + \varepsilon\mathbf{B}^1 + \varepsilon^2\mathbf{B}^2 + \dots) = (\mathbf{B}^0 + \varepsilon\mathbf{B}^1 + \varepsilon^2\mathbf{B}^2 + \dots)(\Lambda_0 + \varepsilon\Lambda_1 + \varepsilon^2\Lambda_2 + \dots). \quad (9)$$

We can look for  $\Lambda_0, \Lambda_1, \Lambda_2, \dots$  by equating the coefficients corresponding to the same powers of  $\varepsilon$ . We start with the system of  $m$  equations what corresponds to the zero power of  $\varepsilon$  in (9)  $\mathbf{A}_0\mathbf{B}^0 - \mathbf{B}^0\Lambda_0 = 0$ . One can satisfy these equations with any basis  $\mathbf{B}^0 = \mathbf{P}(0)\hat{M}(\mathbb{R}^n) \subset \hat{M}(\mathbb{R}^n)$  in the root subspace corresponding to eigenvalue  $r(A_0)^2$  and the matrix  $\Lambda_0$  of the operator  $\mathbf{A}_0$  in this basis. Further we have to deal with the systems of equations which correspond to  $\varepsilon, \varepsilon^2$  and so on in (9). These systems have solutions if and only if its right part is orthogonal to  $m$  linearly independent solutions of homogeneous adjoint equation.

#### 4 Reducibility method for moments

The paper applies reducibility method for moments of equation (1) with near to constant Markov coefficients (7) solutions, which permits to reduce (1) to equation with constant coefficients. It is assumed that a Markov sequence  $\vec{y} := \{y_t, t \in \mathbb{N}\}$  is given in a filtrated probability space  $(\Omega, \mathfrak{F}, \mathfrak{F}^t, P)$ , where  $\{\mathfrak{F}^t\}$  is a minimal filtration adapting it. To write an operator equation for the first moments of (1) in a space of continuous  $n$ -dimensional mappings  $\mathbb{C}(\mathbb{Y} \rightarrow \mathbb{R}^n) := \mathbb{C}_n(\mathbb{Y})$ , a linear continuous operator is introduced:

$$y \in \mathbb{Y}, u \in \mathbb{C}_n(\mathbb{Y}) : (\mathbf{A}u)(y) = \int_{\mathbb{Y}} A^T(z)u(z)p(y, dz). \quad (10)$$

**Lemma 4.1** For any  $s \in \mathbb{R}, t > 0, v \in \mathbb{C}_n(\mathbb{Y}), x \in \mathbb{R}^n$

$$\mathbb{E}\{(X_s^{s+t}x, v(y_{s+t}))/\mathfrak{F}^s\} = (x, (\mathbf{A}^t v)(y_s)).$$

It is said that the equation (1) is mean reducible, if such a continuous matrix function  $\{\mathbf{B}(y), y \in \mathbb{Y}\}$  and such a matrix  $\Lambda$  exist, that for all  $s \in \mathbb{N}$  and  $t > s$  the following equality is fulfilled:  $\mathbb{E}\{\mathbf{B}(y_t)x_t/\mathfrak{F}^s\} = \Lambda^{t-s}\mathbf{B}(y_s)x_s$ .

**Theorem 4.2** Let elements of sequence  $\{y_t, t \in \mathbb{N}\}$  be independent and identically distributed. Then

- (i) operator  $\mathbf{A}$  leaves as invariant a subspace  $\mathbb{R}^n \subset \mathbb{C}_n(\mathbb{Y})$  and restriction  $\bar{\mathbf{A}}$  of operator  $\mathbf{A}$  in this subspace is defined by equality

$$v \in \mathbb{R}^n : \bar{\mathbf{A}}v = \bar{A}^T v,$$

where  $\bar{A} = \mathbb{E}\{A(y_0)\}$ ;



(ii) for each  $s \in \mathbb{N}$ , each  $t > s$  and each  $\mathfrak{F}^t$ -adapted solution  $\{x_t, t \geq 0\}$  of equation (1) the following equality is into force:

$$\mathbb{E}\{x_t\} = \bar{A}^{t-s}\mathbb{E}\{x_s\}.$$

To define the reduced equation of the equation with near to constant matrix coefficients dependent on Markov chain (7) the operator (10) is expressed in a form  $\mathbf{A}(\varepsilon) = \mathbf{A}_0 + \varepsilon\mathbf{A}_1 + \varepsilon^2\mathbf{A}_2 + \dots$ ; hereto, the operator  $\mathbf{A}_0$  leaves as invariant the subspace  $\mathbb{R}^n$ , and it can be represented as a tensor product of operators  $\mathbf{A}_0 = \mathcal{P} \otimes A_0^T$ :

$$h \in \mathbb{C}(\mathbb{Y}), g \in \mathbb{R}^n : \mathbf{A}_0(h \otimes g) = \mathcal{P}h \otimes A_0^T g,$$

where  $\mathcal{P}$  is a Markov operator defined by formula

$$y \in \mathbb{Y}, u \in \mathbb{C}(\mathbb{Y}) : (\mathcal{P}u)(y) = \int_{\mathbb{Y}} u(z)p(y, dz).$$

The tensor representation of the operator allows to simplify the process of finding the spectrum and resolvent, using the spectrum and resolvent of operators which define them. Due to the exponential ergodicity of Markov chain, the spectrum of operator  $\mathbf{A}_0$  can be expressed in a form:

$$\sigma(\mathbf{A}_0) = \{\lambda_1\lambda_2 : \lambda_1 \in \sigma(\mathcal{P}), \lambda_2 \in \sigma(A_0)\} = \sigma(A_0) \cup \sigma_\rho, \quad (11)$$

where  $\sigma_\rho(A_0) := \{\lambda_1\lambda_2 : \lambda_1 \in \sigma(\mathcal{P}), \lambda_2 \in \sigma_\rho\}$ . The main assumption for mean reducibility of the equation (1) is disjunction of sets in spectrum decomposition (11), that is,  $\sigma(A_0) \cap \sigma_\rho = \emptyset$ . It makes possible to offer the asymptotical method which is based on the decomposition of the spectral projection [4] of operator  $\mathbf{A}(\varepsilon)$  by powers of a small parameter  $\varepsilon$ . The conjugated space of  $\mathbb{C}_n(\mathbb{Y})$  is a space of vector-valued measures  $\mathbb{C}_n^*(\mathbb{Y})$ , and the scalar product of elements  $v \in \mathbb{C}_n(\mathbb{Y})$  and  $g \in \mathbb{C}_n^*(\mathbb{Y})$  is defined by the equality  $\langle g, v \rangle := \int_{\mathbb{Y}} (g(dy), v(y))$ .

**Lemma 4.3** *If all the above mentioned assumptions are into force, then for sufficiently small  $\bar{\varepsilon} > 0$  and for all  $|\varepsilon| < \bar{\varepsilon}$ , a difference equation is mean reducible; hereto, the matrix function  $\{\mathbf{B}(y, \varepsilon), y \in \mathbb{Y}\}$  is a basis in operator  $\mathbf{A}(\varepsilon)$  root subspace that corresponds to the part of the spectrum  $\sigma_0(\varepsilon)$  that is defined by equality  $\lim_{\varepsilon \rightarrow 0} \sigma_0(\varepsilon) = \sigma_0$ , but the matrix  $\Lambda(\varepsilon)$  is the operator's  $\mathbf{A}(\varepsilon)$  restriction matrix to this root subspace. For each  $|\varepsilon| < \bar{\varepsilon}$ ,  $n \times n$ -matrix function of the basis  $\{\mathbf{B}(y, \varepsilon), y \in \mathbb{Y}\}$  and constant  $n \times n$ -matrix  $\Lambda(\varepsilon)$  unambiguously are defined by the equality*

$$y \in \mathbb{Y}, |\varepsilon| < \bar{\varepsilon} : (\mathbf{A}(\varepsilon)\mathbf{B})(y, \varepsilon) = \mathbf{B}(y, \varepsilon)\Lambda^T(\varepsilon). \quad (12)$$

The decompositions of the basis matrix  $\mathbf{B}(y, \varepsilon)$  and the matrix  $\Lambda(\varepsilon)$  are used in a form of uniformly converged sequences by powers of a small parameter  $\varepsilon$ :  $\Lambda(\varepsilon) := \Lambda_0 + \varepsilon\Lambda_1 + \varepsilon^2\Lambda_2 + \dots$  and  $\mathbf{B}(y, \varepsilon) := \mathbf{B}_0 + \varepsilon\mathbf{B}_1(y) + \varepsilon^2\mathbf{B}_2(y) + \dots$ . For each sufficiently small  $\varepsilon$  these decompositions can be substituted in the expression (12). Equating the coefficients of equal powers of  $\varepsilon$ , the equations can be obtained for finding the unknown elements of series  $\Lambda(\varepsilon)$  and  $\mathbf{B}(y, \varepsilon)$ :

$$\mathbf{A}_0\mathbf{B}_0 = \mathbf{B}_0\Lambda_0^T \quad (13)$$

$$\mathbf{A}_0\mathbf{B}_1 - \mathbf{B}_1\Lambda_0^T = \mathbf{B}_0\Lambda_1^T - \mathbf{A}_1\mathbf{B}_0 \quad (14)$$

$$\mathbf{A}_0\mathbf{B}_2 - \mathbf{B}_2\Lambda_0^T = \mathbf{B}_0\Lambda_2^T + \mathbf{B}_1\Lambda_1^T - \mathbf{A}_0\mathbf{B}_2 - \mathbf{A}_1\mathbf{B}_1 \quad (15)$$

...

Let us define an operator

$$\begin{aligned} y \in \mathbb{Y}, v \in \hat{\mathbb{C}} : (\mathbb{L}v)(y) &:= (\mathbf{A}_0v)(y) - v(y)A_0^T := \\ &:= \int_{\mathbb{Y}} A_0^T(v(z) - v(y))p(y, dz) + A_0^T v(y) - v(y)A_0^T := \\ &:= (\mathbb{H}v)(y) + (\mathbb{G}v)(y) \end{aligned} \quad (16)$$

for the elements of continuous matrix functions space  $\hat{\mathbb{C}}$ . Looking at  $\hat{\mathbb{C}}$  as at  $\mathbb{R}^{n^2}$ , similarly as in the case with  $\mathbb{C}_n(\mathbb{Y})$ , count additive matrix-valued measure in  $\hat{\mathbb{C}}^*$  can be found, which will be a conjugated space; and the scalar product of elements  $g \in \hat{\mathbb{C}}^*$  and  $v \in \hat{\mathbb{C}}$  can be defined by formula  $\langle g, v \rangle := \text{Tr}\{\int_{\mathbb{Y}} v^T(y)g(dy)\}$ , where  $\text{Tr}\{\}$  is a matrix trace. Taking unit matrix  $\mathbf{B}_0 := I$  as basis in  $\mathbb{R}^n$  and substituting it in the equation (13),  $\Lambda_0^T = A_0^T$  can be found, that is,  $\Lambda_0 = A_0$ . Using Fredholm alternative about normal solvability, the necessary and sufficient conditions can be verified to ensure that a solution exists. The matrix  $\Lambda_2$  can be found from the equation (14), afterwards also  $\mathbf{B}_2(y)$  can be found. Then the next equations can be written for finding  $\Lambda_3, \mathbf{B}_3(y)$  until the necessary accuracy of decomposition of the matrix  $\Lambda(\varepsilon)$  is obtained. Since  $\mathbb{Y}$  is compact and matrices  $\{\mathbf{B}_j(y), j = 1, 2, \dots\}$  are continuous, the elements of the obtained basis  $\mathbf{B} := I + \varepsilon\mathbf{B}_1 + \varepsilon^2\mathbf{B}_2 + \dots$  are linearly independent for sufficiently small  $\varepsilon$ .

## Acknowledgement

The paper was supported from Riga Technical University by project Nr. IZM7388 with title "Copula based auto regressive models of risk prediction" .

## References

- [1] ARNOLD, L., WIHSTUTZ, V. (eds.): *Lyapunov exponents*. Proc., Bremen, 1984. Lecture Notes in Math., vol. 1186, Springer, Berlin, 1986.
- [2] CARKOVA, V., SWERDAN, M.: *On Mean Square Stability of Linear Stochastic Difference Equations*. Theory of Stochastic Ptocesess. TBiMC, vol. 11(27), No 1-2, Kiev. Ukraine, 2005.
- [3] CARKOVA, V., GOLDSTEINE, J.: *Covariance analysis of linear Markov iterations*. System Research and Information Technologies, 3, IASA, Ukraine, 2007.
- [4] KATO, T.: *Perturbations Theory for Linear Operators*. Springer-Verlag, Berlin-Heidelberg, 1966.
- [5] KREIN, M.G., RUTHMAN, M.A.: *The linear operators leaving as invariant cone in Banach space*. Russian Math. Survey 3, No. 1, 3-95, 1947.

**Current address**

**Carkova Viktorija, Dr. Math., Assoc. Prof.**

University of Latvia, 8 Zellu street, Riga, LV-1002, Latvia, tel. number +371 67033727 and e-mail tsarkova@latnet.lv

**Goldšteine Jolanta, Dr. Math.**

University of Latvia, 8 Zellu street, Riga, LV-1002, Latvia, tel. number +371 29498160 and e-mail jolanta.goldsteine@gmail.com

**Swerdan Myhailo, Dr. Phys.-Math. Sci., Prof.**

The State Financial Academy of Bukovina, 1 Shterna street, Chernivtsi, 58000, Ukraine, tel. number +380 372 235152



## ON CONTINUOUS STOCHASTIC MODELING OF HETEROSKEDASTIC CONDITIONAL VARIANCE

CARKOVŠ Jevgenijs, (LV), EGLE Aigars, (LV)

**Abstract.** The proposal continuous stochastic differential equation for conditional variance is constructed as a diffusion approximation of discrete ARCH process. In contrast to classical auto regressive models with independent random perturbations our paper deals with uncertainty given as a stationary ergodic Markov chain. The method is based on stochastic analysis approach to finite dimensional difference equations with proportional to small parameter  $\varepsilon$  increments. Writing a point-form solution of this difference equation as vertexes of a time-dependent continuous broken line given on the segment  $[0,1]$  with  $\varepsilon$ -dependent scaling of intervals between vertexes and tending  $\varepsilon$  to zero we apply probabilistic limit theorems for dynamical systems with rapid Markov switching. The distribution of stationary solution of resulting stochastic equation may be successfully used for analysis of initial discrete model. This method permits to discuss a correlation effect on log of cumulative excess return with stochastic volatility. model-based analysis shows that it is important to take into account possible serial residual correlation in conditional variance process. The proposed method is applied to investigate the GARCH(1,1) and GARCH(2,1) processes under assumption that random variables are serially correlated. As a result it is possible to find continuous stochastic differential equations these processes converge to in distribution.

**Key words and phrases.** Markov dynamical system, diffusion approximation, ARCH model

*Mathematics Subject Classification.* Primary 37H10,37N99; Secondary 34C29,37M10.

### 1 Introduction

Many econometrical studies [1] and [3] have documented that financial time series tend to be highly heteroskedastic. This has many implications for many areas of macroeconomics and finance, including the term structure of interest rates, option pricing and dynamic capital-asset

pricing theory. In the same time econometricians have also been very active in developing models of conditional heteroskedasticity. The most widely used models of dynamic conditional variance have been the ARCH models first introduced by [2]. In most general form, a univariate ARCH model makes conditional variance at time  $t$  a function of exogenous and lagged endogenous variables, time, parameters and past residuals.

In contrast to the stochastic differential equation models so frequently used in theoretical finance literature, ARCH models are discrete time stochastic difference equation systems. Empirics have favored the discrete time approach of ARCH as virtually all time series data are recorded only on discrete time intervals and a discrete time ARCH likelihood function is usually easy to compute and maximize. By contrast, the likelihood of a nonlinear stochastic differential equation system observed at discrete intervals can be very difficult to derive, especially when there are unobservable state variables (like conditional variance) in the system.

Substantial work has been done on relation between the continuous time nonlinear stochastic differential systems, used so much in theoretical literature, and the ARCH stochastic difference equation systems, favored by empirics. Although the two literatures have developed quite independently there have been attempts to reconcile the discrete and continuous models. Nelson [4] is one of the first to partially bridge the gap by developing conditions under which ARCH stochastic difference equation systems converge in distribution to Ito processes as the length of the discrete time intervals goes to zero.

In his work [4] investigates the GARCH(1,1)-M process of [3] for the log of cumulative excess return  $Y_t$ :

$$r_{t+1} = r_t + \varepsilon\gamma\sigma_t^2 + \varepsilon\sigma_t Z_t \quad (1)$$

$$\sigma_{t+1}^2 = \varepsilon^2(\omega - \theta\sigma_t^2) + \varepsilon\alpha\sigma_t^2 Z_t \quad (2)$$

He derives diffusion approximation equation in a form of linear stochastic equation

$$d\sigma_t^2 = (\omega - \theta\sigma_t^2)dt + \alpha\sigma_t^2 dw(t) \quad (3)$$

and shows that in continuous time the stationary distribution for the GARCH(1,1) conditional variance process is an inverted two parametric gamma.

In our paper we change the assumption about independence of random process in equations (1)-(2) and assume that random coefficients in are serially correlated

$$r_{t+1} = r_t + c\sigma_t^2 + \sigma_t y_t \quad (4)$$

$$\sigma_{t+1}^2 = \varepsilon^2(\omega - \theta\sigma_t^2) + \varepsilon\alpha\sigma_t^2 y_t \quad (5)$$

where discrete random process  $y_t$  satisfies AR(1) equation

$$y_t = \rho y_{t-1} + \sqrt{1 - \rho^2} Z_t \quad (6)$$

For analysis of the above equations we apply method and results of the paper [5], which are briefly stated in the first section. In the second section we have derived the stochastic approximation equation for (5) under condition (6) in a form of stochastic equation

$$d\sigma_t^2 = (\omega + (\alpha^2\kappa(\rho) - \theta)\sigma_t^2)dt + \alpha\sqrt{1 - 2\kappa(\rho)}\sigma_t^2 dw(t) \quad (7)$$

with coefficients dependent on correlation parameter  $\rho$  from (6). Analyzing ergodic property and stationary solution of this equation we have shown that it is important to take into account possible serial correlation in conditional variance process. Under the same assumption that random variables are serially correlated we analyze GARCH(2,1) conditional variance process

$$r_{t+1} = r_t + c\sigma_t^2 + \sigma_t y_t \quad (8)$$

$$\sigma_{t+1}^2 = \varepsilon^2(\omega - \theta\sigma_t^2) + \varepsilon\alpha\sigma_t^2 y_t \quad (9)$$

$$y_t = \omega_0 + \alpha_1 y_{t-1} + \alpha_2 y_{t-2} + \beta_1 y_{t-1} Z_t \quad (10)$$

Equations (8) to (10) will be rewritten in a vectorial form and stochastic differential equations will be developed using results found in [5].

## 2 Diffusion Approximation of Discrete Markov Dynamical Systems

### 2.1 Related Work

The problem of asymptotic analysis of dynamical systems under small random perturbations has been discussed in many mathematical and engineering papers. Apparently, A.V. Skorokhod was the first mathematician, which has proved that the probabilistic limit theorems may be successfully used to approximate distributions of solutions of random dynamical systems by the solutions of stochastic differential equations on any finite time interval (see bibliography in [15],[16], and [19]). The above result at once has met with wide application in engineering and economical papers (see [11], [8] [6], [12], [4] and references there). It should be mentioned that in spite of the fact that the above result has been developed for the analysis of equations on a finite time interval, the averaging and diffusion approximation procedures have been applied in many applications for asymptotic stability analysis of possible stationary solutions, that is, for analysis of differential equations as  $t \rightarrow \infty$ . To prove the validity of this approach for random dynamical systems with continuous trajectories the researchers had to use not only a special type of limit theorem (see for example [10] and [7]) but also a stochastic version of the Second Lyapunov method developed for stochastic Ito differential equations in [19]. But most of dynamical systems of the recent Economics (see, for example, [13], [8], [12], [4] and review there) require an extension of the above "smooth" models to allow the phase motion to have a jump type discontinuity. Some of results permissive to resolve this problem have been developed by author in [18] for dynamical systems with switching in Markov time moments. Proposal paper is devoted to similar approach to discrete Markov dynamical systems. This problem is very important in contemporary financial econometrics for analysis of ARCH type stationary iterative procedures (see, for example, [12] and [4]).

### 2.2 Probabilistic limit theorems and equilibrium stochastic stability

Let  $p(y, dz)$  is transition probability of Markov chain  $y_t$  and  $\mathcal{P}$  is Markov operator

$$(\mathcal{P}v)(y) := \int_{\mathbb{Y}} v(z)p(y, dz)$$

defined on the space  $\mathbb{C}(\mathbb{Y})$  of bounded continuous functions. We will assume that the spectrum  $\sigma(\mathcal{P})$  has the simple eigenvalue 1,  $\sigma(\mathcal{P}) \setminus \{1\} \subset \{z \in \mathbb{C} : |z| < \rho < 1\}$ , and probability distribution  $\{\mu(dy)\}$  is the solution of the equation  $\mathcal{P}^*\mu = \mu$ , where  $\mathcal{P}^*$  is conjugate operator. Averaging procedure by the above invariant measure of any dependent on Markov process vector or matrix will be denoted with overline. Under these conditions one can extend [9] the potential of the above Markov process and to define the linear continuous operator by equality

$$(\Pi v)(y) := \sum_{k=0}^{\infty} (\mathcal{P}^k v)(y) \quad (11)$$

on the space  $\bar{\mathbb{C}}(\mathbb{Y})$  of continuous functions  $v \in \mathbb{C}(\mathbb{Y})$  with zero average  $\bar{v} := \int_{\mathbb{Y}} v(y) \mu(dy)$ . This means that the equation  $\mathcal{P}g - g = -v$  with  $v \in \bar{\mathbb{C}}(\mathbb{Y})$  has unique solution (11) in  $\bar{\mathbb{C}}(\mathbb{Y})$ . Using the above Markov chain one can define on the segment  $[0, 1]$  step processes

$$s \in [t\varepsilon^2, (t+1)\varepsilon^2) : Y_\varepsilon(s) := y_t \quad (12)$$

If  $\mathfrak{F}^t \subset \mathfrak{F}$ ,  $t \geq 0$  is minimal filtration for stationary process  $y_t$  then for any  $t \geq 0$  and  $s \in [t\varepsilon^2, (t+1)\varepsilon^2)$  random vectors  $X_\varepsilon(s)$  and  $Y_\varepsilon(s)$  are  $\mathfrak{F}^t$ -measurable. To avoid cumbersome formulae we will denote conditional expectation  $\mathbf{E}\{\xi/\mathfrak{F}^t\}|_{x_t=x, y_t=y}$  in abridged form  $\mathbf{E}_{x,y}^t\{\xi\}$ .

### 2.3 Derivation of diffusion approximation formula

In this subsection we will assume that  $\bar{f}_1(x) \equiv 0$ . Using the solution  $x_t, t \in \mathbb{N}$  of difference equation

$$x_{t+1} = x_t + \varepsilon f_1(x_t, y_t) + \varepsilon^2 f_2(x_t, y_t), \quad (13)$$

with initial condition  $x_0 = x$  and Markov process  $y_t$  one can define the broken lines by formulae

$$s \in [t\varepsilon^2, (t+1)\varepsilon^2) : X_\varepsilon(s) = (x_{t+1} - x_t)(s\varepsilon^{-2} - t) + x_t \quad (14)$$

for all  $t \in [0, N(\varepsilon^{-2})]$ , where  $N(\alpha)$  is integer part of number  $\alpha$ , and step process

$$s \in [t\varepsilon^2, (t+1)\varepsilon^2) : Y_\varepsilon(s) := y_t \quad (15)$$

for all  $t \in [0, N(\varepsilon^{-2})]$ . Not so difficult to be certain of Markov properties for the pair  $\{X_\varepsilon(s), Y_\varepsilon(s), 0 \leq s \leq 1\}$ . Therefore under assumption that  $\varepsilon \rightarrow 0$  one can apply the Skorokhod limit theorems from [15] and [16] for sequences of Markov processes and look for diffusion approximation of  $\{X_\varepsilon(s), 0 \leq s \leq 1\}$  if the latter exists. Much as it has been done in [18] for jump type Markov processes in our case for any arbitrary twice continuous differentiable on  $x$  function  $v(x)$  one has to look for Lyapunov function in a form of decomposition

$$v^\varepsilon(x, y) := v(x) + \varepsilon[(\Pi f_1)(x, y), \nabla]v(x, y) + \varepsilon^2 \hat{v}(x, y) \quad (16)$$

with some smooth function  $\hat{v}(x, y)$ . Here and further  $\nabla v(x)$  is gradient and  $(\cdot, \cdot)$  is scalar product in  $\mathbb{R}^d$ . Now one should compute derivative

$$\begin{aligned} (\mathbf{L}(\varepsilon)v^\varepsilon)(x, y) &:= \\ \lim_{\delta \downarrow 0} \frac{1}{\delta} \mathbf{E}_{x,y}^t \{v^\varepsilon(X^\varepsilon(s+\delta), Y^\varepsilon(s+\delta)) - v^\varepsilon(X^\varepsilon(s), Y^\varepsilon(s))\} &= \\ \frac{1}{\varepsilon^2} \mathbf{E}_{x,y}^t \{v^\varepsilon(x_{t+1}, y_{t+1}) - v^\varepsilon(x, y) + o(\varepsilon^2)\} & \end{aligned} \quad (17)$$



for all  $x \in \mathbb{R}^d, y \in \mathbb{Y}, t \geq 0$  and  $s \in [t\varepsilon^2, (t+1)\varepsilon^2)$ , and chose in (16) function  $\hat{v}(x, y)$  in such a way as to exist limit

$$\lim_{\varepsilon \rightarrow 0} (\mathbf{L}(\varepsilon)v^\varepsilon)(x, y) = (\mathbf{L}v)(x) \quad (18)$$

As it will be shown later right side of the above equation has a form of diffusion operator applied to function  $v(x)$ :

$$(\mathbf{L}v)(x) = \{(a(x), \nabla) + (\sigma(x)\nabla, \nabla)\}v(x) \quad (19)$$

with vector  $a(x)$  and positive defined symmetric matrix  $\sigma(x)$ . To derive the above formula one has to present operator  $\mathbf{L}(\varepsilon)$  accurate within  $0(\varepsilon)$

$$\begin{aligned} \mathbf{L}(\varepsilon) &= \frac{1}{\varepsilon^2}(\mathcal{P} - I) + \frac{1}{\varepsilon}(f_1(x, y), \nabla)\mathcal{P} + \\ & (f_2(x, y), \nabla)\mathcal{P} + \frac{1}{2}(f_1(x, y), \nabla)^2\mathcal{P} + 0(\varepsilon) \end{aligned} \quad (20)$$

to employ (20) to (16) and to decompose resulting function by powers of  $\varepsilon$  accurate within  $0(\varepsilon)$ :

$$\begin{aligned} (\mathbf{L}(\varepsilon)v^\varepsilon)(x, y) &= \frac{1}{\varepsilon^2}(\mathcal{P} - I)v(x) + \\ & \frac{1}{\varepsilon}[(f_1(x, y), \nabla)v(x) + (\mathcal{P} - I)((\Pi f_1)(x, y), \nabla)v(x)] + \\ & (f_2(x, y), \nabla)v(x) + \frac{1}{2}(f_1(x, y), \nabla)^2v(x) + \\ & (f_1(x, y), \nabla)\mathcal{P}[(f_1(x, y), \nabla)v(x)] + (\mathcal{P} - I)\hat{v}(x, y) + 0(\varepsilon) \end{aligned}$$

Therefore using obvious equalities  $(\mathcal{P} - I)\Pi = -I$ ,  $(\mathcal{P} - I)v(x) = 0$  and formula (18) one can write equation

$$\begin{aligned} \mathbf{L}v(x) &= (f_2(x, y), \nabla)v(x) + \frac{1}{2}(f_1(x, y), \nabla)^2v(x) + \\ & (f_1(x, y), \nabla)[(\mathcal{P}\Pi f_1(x, y), \nabla)v(x)] + (\mathcal{P} - I)\hat{v}(x, y) \end{aligned}$$

with unknown function  $\hat{v}(x, y)$ . As it has been mentioned at the beginning of this subsection the above equation relative to  $\hat{v}(x, y)$  has solution

$$\begin{aligned} \hat{v}(x, y) &= \Pi\{(f_2(x, y), \nabla)v(x) + \frac{1}{2}(f_1(x, y), \nabla)^2v(x) + \\ & (f_1(x, y), \nabla)[(\mathcal{P}\Pi f_1(x, y), \nabla)v(x)] - \mathbf{L}v(x)\} \end{aligned} \quad (21)$$

if and only if

$$\mathbf{L}v(x) = \{(\overline{f_2}, \nabla) + \frac{1}{2}(\overline{f_1}, \nabla)^2 + \overline{(f_1, \nabla)[(\mathcal{P}\Pi f_1, \nabla)]}\}v(x) \quad (22)$$

where overline denotes averaging by measure  $\mu$ . This equation one can write in a form (19) using notations

$$\begin{aligned} a &= \overline{f_2} + \overline{[\mathcal{P}\Pi D f_1]^T f_1} \\ \sigma &= \frac{1}{2}[\overline{f_1 f_1^T} + \overline{f_1 \mathcal{P}\Pi f_1^T} + \overline{(\mathcal{P}\Pi f_1) f_1^T}] \end{aligned} \quad (23)$$

where  $D$  is Frechet derivative by  $x$  and upper index  $T$  denotes transposition. To write this equation in a form

$$dX(t) = a(X(s))ds + \sum_{k=1}^d \sigma_k(X(s))dW_k(s) \quad (24)$$

with initial condition  $X(0) = x_0$ , where vector-functions  $a(x)$  and  $\sigma_k(x)$ ,  $k = 1, 2, \dots, d$  are defined based on averaging by measure  $\mu$  of functions  $f_j(x, y)$ ,  $j = 1, 2$  and its derivatives, and  $\{W_k, k = 1, 2, \dots, d\}$  are independent standard Wiener processes, one has to find  $d$  dependent on  $x$  vectors  $\sigma_k$  defined by equation

$$\sum_{k=1}^d \sigma_k(x) \sigma_k^T(x) = \sigma(x)$$

As it has been mentioned in [19] this equation has solution for any positive defined matrix  $\sigma(x)$ .

## 2.4 Averaging and normalized deviations

Let us remind of assumption  $\bar{f}_1(x) \equiv 0$  which permits in previous subsection to derive formulae (22) and (23). Otherwise one may not divide segment  $[0, 1]$  by intervals of length  $\varepsilon^2$  because  $\Pi f_1(x, y)$  does not exist and therefore there are singularity in the definition of operator (17) as  $\varepsilon \rightarrow 0$ . To apply a diffusion approximation in this case one has to find solution of averaged equation

$$\bar{x}_{t+1} = \bar{x}_t + \varepsilon \bar{f}_1(\bar{x}_t) \quad (25)$$

and to derive an asymptotic formula for so called *normalized deviations*

$$z_t := \frac{x_t - \bar{x}_t}{\sqrt{\varepsilon}} \quad (26)$$

Substituting  $x_t = \sqrt{\varepsilon} z_t + \bar{x}_t$  in (13)

$$z_{t+1} = z_t + \delta g_1(\bar{x}_t, y_t) + \delta^2 [Df_1(\bar{x}_t, y_t)] z_t + o(\delta^2), \quad (27)$$

where  $\delta = \sqrt{\varepsilon}$ ,  $g_1(x, y) = f_1(x, y) - \bar{f}_1(x)$ , one can apply to system (25)-(27) approach of previous subsection. The sequence (26) gives rise to random processes

$$Z^\delta(s) = \frac{X^\delta(s) - \bar{X}^\delta(s)}{\delta}$$

where  $X^\delta(s)$  and  $\bar{X}^\delta(s)$  are defined in the same way like (14) for all  $s \in [t\delta^2, (t+1)\delta^2]$  and any  $t \in [0, N(\delta^{-2})]$ . After substitution  $Z^\delta(s)$  instead of  $X_\varepsilon(s)$ ,  $[Df_1(\bar{X}(s), Y^\delta(s))]Z^\delta(s)$  instead of  $f_2(X_\varepsilon(s), Y_\varepsilon(s))$  and  $g_1(\bar{X}(s), Y^\delta(s))$  instead of  $f_1(X_\varepsilon(s), Y_\varepsilon(s))$  in corresponding formulae and vanishing  $\delta$  one can approximate probability distribution  $\mathbf{P}_\delta^Z$  of process  $Z^\delta(s)$  by probability distribution  $\mathbf{P}^Z$  of process  $Z$  satisfying stochastic differential equation

$$dZ(s) = D\bar{f}_1(\bar{X}(s))Z(s)ds + \sum_{k=1}^d \sigma_k(\bar{X}(s))dW_k(s)$$

with initial condition  $\bar{X}(0) = x_0$ , where  $\{W_k(s), k = 1, 2, \dots, d\}$  are independent standard Wiener processes, and vectors  $\{\sigma_k, k = 1, 2, \dots, d\}$  satisfy an equality

$$\sum_{k=1}^d \sigma_k(x) \sigma_k^T(x) = [\overline{g_1 g_1^T} + \overline{g_1 \mathcal{P} \Pi g_1^T} + \overline{(\mathcal{P} \Pi g_1) g_1^T}](x)$$

Deterministic function  $\bar{X}(s)$  one can find as the solution of ordinary differential equation

$$d\bar{X}(s) = \bar{f}_1(\bar{X}(s))ds$$

Roughly speaking for sufficiently small  $\varepsilon$  one can approximate distribution of the sequence  $\{x_t, 0 \leq t \leq N(\varepsilon^{-1})\}$  by distribution of sequence  $\{X(t\varepsilon) + \sqrt{\varepsilon}Z(t\varepsilon), 0 \leq t \leq N(\varepsilon^{-1})\}$ .

## 2.5 Equilibrium asymptotic stability

As it has been mentioned in the Section 2 some of application iterative procedures analysis require asymptotic analysis of equation (13) as  $t \rightarrow \infty$ . For example discussing diffusion approximation approach to GARCH time series authors of papers [4] and [12] indicate this problem in view of the approximation and asymptotic stability analysis of stationary conditional variance. In previous section we have derived an approximate distribution of sequence  $\{x_t, 0 \leq t \leq N\}$  for any finite integer number  $N$  by distribution of solution of stochastic differential equation  $\{X(s), 0 \leq s \leq 1\}$  but for the above mentioned asymptotic analysis as  $t \rightarrow \infty$  one has to deal with equation (24) with unrestrictedly large  $s$ . Besides there is a problem of legality results which are based on the diffusion approximation as  $s \rightarrow \infty$ . This subsection is devoted to the above problem.

Let point  $x = 0$  be an equilibrium of iteration procedure (13), i.s.  $f_1(0, y) \equiv 0$  and  $f_2(0, y) \equiv 0$ . If for any  $\eta > 0$  there exists such a neighborhood  $U_\eta := \{x \in \mathbb{R}^d : |x| < \eta\}$  that any starting in  $U_\eta$  solution  $x_t$  of (13) does not leave  $U_\eta$  and tends to zero as  $t \rightarrow \infty$  with probability greater than  $1 - \eta$  the above equilibrium is called *asymptotic stochastically stable*. As it has been shown in [14] for equilibrium stability analysis one can employ the second Lyapunov method with Lyapunov operator defined by formula

$$(\mathcal{L}v)(x, y) := \mathbf{E}_{x,y}^0 \{v(x_1, y_1)\} - v(x, y)$$

and Lyapunov functions satisfying inequality

$$|x|^p < v(x, y) < c|x|^p$$

with some positive  $p$  a  $c \geq 1$ . If there exists such a Lyapunov function  $v(x, y)$  that

$$(\mathcal{L}v)(x, y) < -\gamma|x|^p$$

with  $\gamma \in (0, 1)$  then [14] equilibrium is asymptotic stochastically stable and  $\mathbf{E}_{x,y} \{|x_t|\} \leq M|x|^p \exp\{-\rho t\}$  with some positive constants  $M$  and  $\rho$ . Besides under smoothness assumptions of the Section 1 on vectors  $f_1(x, y)$  and  $f_2(x, y)$ , this equilibrium is asymptotic stochastically stable if and only if [14] the same property has the trivial solution of its linear approximation

$$\tilde{x}_{t+1} = \tilde{x}_t + \varepsilon \tilde{f}_1(\tilde{x}_t, y_t) + \varepsilon^2 \tilde{f}_2(\tilde{x}_t, y_t) \quad (28)$$

where  $\tilde{f}_j(x, y) = (Df_j)(0, y)x$ ,  $j = 1, 2$ . Therefore for asymptotic analysis of (13) as  $t \rightarrow \infty$  one can apply formulae (19) with (16), (21), (22), and (23) substituting linear on  $x \in \mathbb{R}^d$  functions  $\tilde{f}_j(x, y)$  instead of  $f_j(x, y)$ ,  $j = 1, 2$  and rewriting equation (24) in a form of linear stochastic Ito equation

$$d\tilde{X}(s) = A\tilde{X}(s)ds + \sum_{k=1}^d B_k\tilde{X}(s)dW_k(s) \quad (29)$$

The same result like mentioned above for Markov iterations (28) one can find in [19] for stochastic differential equation (29): trivial solution of (29) is asymptotic stochastically stable if and only if there exists such twice continuous differentiable Lyapunov function  $V(x)$  that

$$|x|^p \leq V(x) \leq h_1|x|^p, \quad \mathbf{L}V(x) \leq -h_2|x|^p \quad (30)$$

and  $\|D^l \nabla v(x)\| \leq h_3|x|^{p-l-1}$ ,  $l = 1, 2, 3$  for some  $p > 0$ , positive constants  $h_j$ ,  $j = 1, 2, 3$ . and any  $x \in \mathbb{R}^d$ . Now for analysis of asymptotic behaviour of linear iteration (28) one can apply the second Lyapunov method with function

$$V^\varepsilon(x, y) := V(x) + \varepsilon[(\Pi\tilde{f}_1)(x, y), \nabla]V(x, y) + \varepsilon^2\hat{V}(x, y) \quad (31)$$

where  $V(x)$  satisfies inequalities (30) and

$$\begin{aligned} \hat{V}(x, y) = & \Pi\{(\tilde{f}_2(x, y), \nabla)V(x) + \frac{1}{2}(\tilde{f}_1(x, y), \nabla)^2V(x) + \\ & (\tilde{f}_1(x, y), \nabla)[(\mathcal{P}\Pi\tilde{f}_1)(x, y), \nabla]V(x)\} \end{aligned} \quad (32)$$

Owing to linearity of functions  $\tilde{f}_j(x, y)$ ,  $j = 1, 2$  and definition of  $\mathbf{L}V(x)$  for all sufficiently small  $\varepsilon > 0$  there exist such positive constants  $h_j$ ,  $j = \overline{4, 9}$  that the above defined functions satisfy inequalities

$$\begin{aligned} h_4|x|^p & \leq |\hat{V}(x, y)| \leq h_5|x|^p, \\ h_6|x|^p & \leq |[(\Pi\tilde{f}_1)(x, y), \nabla]V(x, y)| \leq h_7|x|^p \\ |V^\varepsilon(x, y) - V(x)| & \leq \varepsilon h_8|x|^p \end{aligned}$$

and

$$|(\mathbf{L}(\varepsilon)V^\varepsilon)(x, y) - \mathbf{L}V(x)| < \varepsilon h_9|x|^p$$

Therefore if the trivial solution of diffusion approximation is asymptotically stable then there exists Lyapunov function satisfying (30) and for stability analysis of (28) one can use function (31):

$$\begin{aligned} (\mathcal{L}V^\varepsilon)(x, y) & = \varepsilon^2(\mathbf{L}(\varepsilon)V^\varepsilon)(x, y) \leq \\ & \leq \varepsilon^2\mathbf{L}V(x) \leq \varepsilon^2(-h_2 + \varepsilon h_9)|x|^p \end{aligned}$$

This inequality convinces of asymptotical stochastic stability for trivial solution of difference equation (28) if  $\varepsilon$  is sufficiently small.

### 3 Markov type GARCH models

#### 3.1 Continuous Stochastic Model of Conditional Variance Dynamics

In papers [12] and [4] the authors discuss a problem of diffusion approximation for very popular in contemporary econometrics GARCH (General AutoRegressive Conditional Heteroscedastic) process for conditional time series variance. The paper [4] deals with model given in a form of first order linear difference equation

$$\sigma_{t+1}^2 = \omega_h + \sigma_t^2[\beta_h + h^{-1}\alpha_h Z_t^2] \quad (33)$$

where  $h$  is small positive parameter,  $\{Z_t, t \in \mathbb{Z}\}$  is sequence of i.i.d. random variables with zero mean, variance  $\mathbf{E}\{Z_t^2\} = h$ , and fourth moment  $\mathbf{E}\{Z_t^4\} = 3h^2$ . Under assumptions

$$1 - \alpha_h - \beta_h = h\theta + o(h), \omega_h = h\omega + o(h), \alpha_h = \frac{\sqrt{h}}{\sqrt{2}}\alpha + o(h)$$

author of paper [4] derives diffusion approximation equation in a form

$$d\sigma_t^2 = (\omega - \theta\sigma_t^2)dt + \alpha\sigma_t^2 dW(t) \quad (34)$$

To compare this result with our derived formulae one can denote

$$h = \varepsilon^2, x_t = \sigma_t^2, y_t = \frac{Z_t^2 - 1}{\sqrt{2h}}$$

and rewrite equation (33) in a form of difference equation (13) accurate within  $\varepsilon$ -items of second order

$$x_{t+1} = x_t + \varepsilon^2[\omega - \theta x_t] + \varepsilon\alpha y_t x_t \quad (35)$$

Let  $y_t$  be stationary Markov process with the same unconditional moments as  $\frac{Z_t^2 - 1}{\sqrt{2h}}$ , that is,  $\mathbf{E}y_t = 0$ ,  $\mathbf{E}y_t^2 = 1$  and correlation function  $C(k) = \mathbf{E}\{y_t y_{t+k}\}$  for  $k \in \mathbb{N}$ . Following our proposal method of diffusion approximation one should for this equation calculate parameters (23) with  $f_1(x, y) = \alpha y x$ ,  $f_2(x, y) = \omega - \theta x$ . By definition

$$\begin{aligned} a(x) &= \omega - \theta x + \alpha^2 x \sum_{l=1}^{\infty} \left\{ \int_{\mathbb{Y}} \mathbf{E}_y^0\{y y_l\} \mu(dy) \right\} = \\ &= \omega + \left[ \alpha^2 \sum_{k=1}^{\infty} C(k) - \theta \right] x \end{aligned}$$

$$\begin{aligned} \sigma^2(x) &= \alpha^2 x^2 \int_{\mathbb{Y}} y^2 \mu(dy) + 2\alpha^2 x^2 \sum_{k=1}^{\infty} C(k) = \\ &= \alpha^2 x^2 \left[ \text{Var}\{y_t\} + 2 \sum_{k=1}^{\infty} C(k) \right] \end{aligned}$$

If  $\{y_t, t \in \mathbb{Z}\}$  are independent random variables with zero mean and unit variance like it has been assumed in [4] we have derived equation (34) because  $C(k) \equiv 0$ . If  $\kappa := \sum_{k=1}^{\infty} C(k) \neq 0$  one should apply diffusion approximation for GARCH(1,1)-process in a following form

$$d\sigma_t^2 = (\omega + (\alpha^2\kappa - \theta)\sigma_t^2)dt + \alpha\sqrt{1 + 2\kappa}\sigma_t^2 dW(t) \quad (36)$$

As it has been proved this equation one can use also for analysis of (33) as  $t \rightarrow \infty$ . According to [19] if

$$\alpha^2\kappa - \theta - \frac{\alpha^2(1 + 2\kappa)}{2} = -\theta - \frac{\alpha^2}{2} < 0$$

there exists stationary solution  $\hat{s}_t^2$  of the above equation and deviations  $z_t := s_t^2 - \hat{s}_t^2$  of any other solution from this stationary process exponentially tend to zero as  $t \rightarrow \infty$ . In spite of the fact that process  $y_t$  has nonzero correlation this result no differs from similar result of the paper [4]. But to approximate stationary process for GARCH(1,1) with Markov process  $y_t$  instead of i.i.d. sequence one has to deal with stationary solution of equation (36) where  $\kappa \neq 0$ . As it has been derived by E.Wong [17] for linear stochastic Ito equation (36) the stationary process defined  $\sigma_t^{-2}$  has density function  $f(x) = \frac{s^r x^{(r-1)}}{\Gamma(r)} e^{-sx}$  where  $r = 1 + \frac{2(\theta - \alpha^2\kappa)}{\alpha^2(1+2\kappa)}$ ,  $s = \frac{2\omega}{\alpha^2(1+2\kappa)}$  or stationary variance  $\hat{s}^2$  has distribution defined by formula

$$\mathbb{P}\{\hat{s}^2 < z\} = 1 - \int_0^{1/z} f(x)dx, \quad f(x) = \frac{s^r x^{(r-1)}}{\Gamma(r)} e^{-sx} \quad (37)$$

This convince of possible considerable correlation affect on the asymptotic approximation of conditional variance stationary distribution.

### 3.2 Diffusion Model of Stock Return with Stochastic Volatility

The simplest stock return  $S_t$  mathematical model involving assumption on conditional heteroskedasticity of interest rate  $h_t$  variance  $\sigma_t^2$  under commonly used condition on risk neutrality of probabilistic measure  $\mathbf{P}$  may be written ([2],[4]) as the system of two difference equation

$$S_{t+1} = S_t(1 + \varepsilon\sigma_t^2 y_{t+1}), \quad (38)$$

$$\sigma_{t+1}^2 = \sigma_t^2 + \varepsilon^2[\omega - \theta\sigma_t^2] + \varepsilon\alpha(y_{t+1}^2 - 1)\sigma_t^2 \quad (39)$$

where  $y_t$  is Gaussian random sequence with zero mean and unit variance. When it is considered that these random numbers do not independent we will use for  $y_t$  equation of type AR(1):

$$y_{t+1} = \rho y_t + \sqrt{1 - \rho^2} \xi_{t+1} \quad (40)$$

where  $\{\xi_t\}$  is i.i.d. Gaussian sequence,  $\mathbb{E}\xi_t = 0, \mathbb{E}\xi_t^2 = 1$ . To employ formulae (refLL) let us denote  $x_{1t} = S_t, x_{2t} = \sigma_t^2$  and

$$\vec{x}_t = \begin{pmatrix} x_{1t} \\ x_{2t} \end{pmatrix}$$

and rewrite equations (38) in a vector form

$$\vec{x}_{t+1} = \vec{x}_t + \varepsilon y_{t+1} \begin{pmatrix} \sqrt{x_{2t}} & 0 \\ 0 & \alpha \end{pmatrix} \vec{x}_t + \varepsilon^2 \begin{pmatrix} 0 \\ \omega \end{pmatrix} - \varepsilon^2 \begin{pmatrix} 0 & 0 \\ 0 & \theta \end{pmatrix} \vec{x}_t \quad (41)$$

Now one can use formula (refLL) with

$$f_1 = \begin{pmatrix} x_{1t}y_{1t}\sqrt{x_{2t}} \\ \alpha x_{2t}(y_t^2 - 1) \end{pmatrix}, \quad f_2 = \begin{pmatrix} 0 \\ \omega - \theta x_{2t} \end{pmatrix}$$

applying in formulae (refLL) an averaging by invariant distribution  $N(0, 1)$  of Markov chain (40) and to write a final limit stochastic equation for vector  $\vec{x}_t$ :

$$d\vec{x}_t = a(\vec{x}_t)dt + b_1(\vec{x}_t)dw_1(t) + b_2(\vec{x}_t)dw_2(t)$$

where

$$a(\vec{x}) = \begin{pmatrix} x_1x_2\frac{\rho}{1-\rho} \\ \omega + (\alpha^2\frac{\rho^2}{1-\rho^2} - \theta)x_2 + \frac{x_1^2}{2} \end{pmatrix} \quad (42)$$

vectors  $b_1, b_2$  are defined by equality

$$b_1(\vec{x})b_1^T(\vec{x}) + b_2(\vec{x})b_2^T(\vec{x}) = \begin{pmatrix} x_1^2x_2\frac{3+\rho}{1-\rho} & 0 \\ 0 & 2\alpha^2x_2^2\frac{3+\rho^2}{1-\rho^2} \end{pmatrix} \quad (43)$$

This means that stock return  $S_t$  and conditional variance  $\sigma_t^2$  have dynamics approximately describes by system of Ito stochastic differential equations

$$dS_t = S_t\sigma_t\frac{\rho}{1-\rho}dt + S_t\sigma_t\sqrt{\frac{3+\rho}{1-\rho}}dw_1(t), \quad (44)$$

$$d\sigma_t^2 = \{\omega + (\alpha^2\frac{\rho^2}{1-\rho^2} - \theta)\sigma_t^2 + S_t^2\frac{\rho}{2(1-\rho)}\}dt + \alpha\sigma_t^2\sqrt{\frac{2(3+\rho^2)}{1-\rho^2}}dw_2(t) \quad (45)$$

where  $w_1(t)$  and  $w_2(t)$  are independent standard Wiener processes.

## Acknowledgement

The paper was supported by grant from Latvian Scientific Council no. 05.1879 with title "Asymptotic analysis of stochastic stability" and by project from Riga Technical University no. 7388 with title "Copula based auto regressive models of risk prediction" .

## References

- [1] BLACK, F.: *Studies of Stock market Volatility Changes*. In Proceedings of the American Statistical Association, Business and Economic Statistics Section, pp. 177-181, 1976.
- [2] ENGLE, R.F.: *Autoregressive Conditional Heteroskedasticity with Estimates of the Variance of United Kingdom Inflation*. In *Econometrica*, 50, pp. 987-1008, 1982.

- [3] ENGLE, R.F. and BOLRSLEV, T.: *Modelling the Persistence of Conditional Variances*. In Economic Reviews, 5, pp. 1-50, 1986
- [4] NELSON, D.B. *ARCH Models as Diffusion Approximation*, In Journ. of Econometrics. 45, pp. 7-38, 1990
- [5] CARKOV, J.: *On Diffusion Approximation of Discrete Markov Dynamical Systems*. In Computational Geometry, Proceedings of World Academy of Science, Engineering and Technology, Vol. 30, July, pp. 1-6, 2008.
- [6] ANN, C.M. and THOMPSON, H.: *Jump-diffusion process and term structure of interest rates*. In J. of Finance, 43, No. 3, pp. 155-174, 1988.
- [7] ARNOLD, PAPANICOLAOU, L. G., and WIHSTUTZ, V.: *Asymptotic analysis of the Lyapunov exponent and rotation number of the random oscillator and applications*, In SIAM J. Appl. Math., 46, No. 3, pp. 427-450, 1986.
- [8] BALL, C.A. and ROMA, A.C.: *A jump diffusion model for the European Monetary System*. In J. of International Money and Finance, pp. 475-492, 1993.
- [9] DYNKIN, E.B.: *Markov Processes*, NY, USA: Academic Press, 1965.
- [10] BLANKENSHIP, G. and PAPANIKOLAOU, G.: *Stability and control of stochastic systems with wide-band noise disturbances.I*. In SIAM J. Appl. Math, 34, No. 3, pp. 437-476, 1978.
- [11] DIMENTBERG, M.: *Statistical Dynamics of Nonlinear and Time-Varying Systems*. NY, USA: Willey, 1988.
- [12] FORNARI, F. and MELE, A.: *Sign- and volatility-switching ARCH models theory and applications to international stock markets*. In J. of Applied Econometrics, 12, pp. 49-65, 1997.
- [13] JARROW, R. and ROSENFELD, E.: *Jump risks and the International capital asset pricing models*. In J. of Business, 57, pp. 337-351, 1984.
- [14] M.B. Nevelson and R.Z. Hasminskii. *Stochastic Approximation and Recursive Estimation*, NY, USA: American Mathematical Society, ISBN-10:0821809067, 1976
- [15] SKOROKHOD, A.V.: *Studies in the theory of random processes*, , USA: Addison-Wesley publishing, 1965.
- [16] SKOROKHOD, A.V.: *Asymptotic Methods of Theory of Stochastic Differential Equations*, 3rd ed. AMS, USA: Providance, 1994.
- [17] WONG, E.: *The construction of a class of stationary Markov process* In (R.Bellman ed.), Proceedings of the 16th Symposia in Applied Mathematics: Stochastic Processes in Mathematical Physics and Engineering, Providance, RI., pp. 264-276, 1964.
- [18] TSARKOV, Ye. (J.Carkovs): *Asymptotic methods for stability analysis of Markov impulse dynamical systems*, In Nonlinear dynamics and system theory, 1, No. 2, pp. 103 - 115, 2002.
- [19] KHASHINSKY, R.Z.: *Limit theorem for a solution of the differential equation with a random right part*. In Prob. Theor. and its Appl., 11, No 3, pp. 444 - 462, 1966.



**Current address**

**Jevgenijs Carkovs, prof.**

Riga Technical University, Kalķu 1, LV-1050, Riga, LATVIA, phone: +371 67089517,  
e-mail: carkovs@latnet.lv

**Aigars Egle, PhD st.**

Riga Technical University, Kalķu 1, LV-1050, Riga, LATVIA, phone: +371 67089517,  
e-mail: carkovs@latnet.lv



# APPLICATION OF INTEGRAL INEQUALITIES IN THE THEORY OF INTEGRAL AND INTEGRODIFFERENTIAL EQUATIONS

FAJMON Břetislav, (CZ), ŠMARDA Zdeněk, (CZ)

**Abstract.** In this paper boundedness and asymptotic stability of solutions of certain classes of integral and integrodifferential equations are investigated. By means of inequalities with iterated integral there are determined conditions of boundedness and asymptotic stability of solutions of nonlinear Volterra integral equations and using of Pachpatte's integral inequalities sufficient conditions of boundedness of solutions of certain class of nonlinear integrodifferential equations are given.

**Key words and phrases.** Inequalities with iterated integrals, boundedness and asymptotic stability of solutions, integral and integrodifferential equations.

*Mathematics Subject Classification.* 45J05.

## 1 Introduction

Integral inequalities involving functions and their derivatives have played a significant role in the developments of various branches of analysis. Pachpatte[6-10] has given some integral inequalities of the Gronwall-Bellman type involving functions and their derivatives which are useful in the investigation of boundedness and stability of solutions of differential, integral and integrodifferential equations. Haraux[4] used the modified Gronwall-Bellmann inequality with logarithmic factor in the integrand to the study of wave equation with logarithmic nonlinearity. Engler[2] used the slight variant of the Haraux's inequality for determination of global regular solutions of the dynamic antiplane shear problem in nonlinear viscoelasticity. Dragomir[1] applied his inequality to the stability, boundedness and asymptotic behaviour of solutions of nonlinear Volterra integral equations.

In this paper we present applications of some above mentioned inequalities to certain integral and integrodifferential equations.

## 2 Bounds of solutions of integrodifferential equations

Consider the following integrodifferential equation

$$x'(t) - F\left(t, x(t), \int_0^t k(t, s, x(s))ds\right) = h(t), \quad x(0) = x_0, \quad (1)$$

where  $h : R^+ \rightarrow R$ ,  $k : R^+ \times R^+ \times R \rightarrow R$ ,  $F : R^+ \times R^2 \rightarrow R$  are continuous functions. Some classes of equation (1) were investigated by Yang[11].

In the following we will suppose that the solution  $x(t)$  of (1) exists on  $R^+$ .

Now for the determination of a bound of the solution  $x(t)$  of (1) we utilize Pachpatte's inequality [10] which we can formulate as the following theorem.

**Theorem 2.1** *Let  $u$ ,  $f$ ,  $g$ ,  $h$  be nonnegative continuous functions defined on  $R^+$  and  $c$  be a nonnegative constant.*

(I) *If*

$$u^2(t) \leq c^2 + 2 \int_0^t \left[ f(s)u(s) \left( u(s) + \int_0^s g(\tau)u(\tau)d\tau \right) + h(s)u(s) \right] ds, \quad (2)$$

*for  $t \in R^+$ , then*

$$u(t) \leq p(t) \left[ 1 + \int_0^t f(s) \exp \left( \int_0^s [f(\tau) + g(\tau)]d\tau \right) ds \right]$$

*where*

$$p(t) = c + \int_0^t h(s)ds \quad (3)$$

*for  $t \in R^+$ .*

(II) *If*

$$u^2(t) \leq c^2 + 2 \int_0^t \left[ f(s)u(s) \left( \int_0^s g(\tau)u(\tau)d\tau \right) + h(s)u(s) \right] ds,$$

*for  $t \in R^+$ , then*

$$u(t) \leq p(t) \exp \left( \int_0^t f(s) \left( \int_0^s g(\tau)d\tau \right) ds \right),$$

*for  $t \in R^+$ , where  $p(t)$  is defined by (3).*

Pachpatte [8] applied these inequalities to obtain bounds of solutions of some classes of ordinary differential equations in the form

$$x'' = G(t, x(t), x'(t)),$$

with initial conditions

$$x(t_0) = x_0, \quad x'(t_0) = x_1,$$

where  $G : I \times R \times R \rightarrow R$  is a continuous function and  $x_0, x_1$  are constants,  $I = [t_0, \infty)$ ,  $t_0 \geq 0$ .

**Theorem 2.2** Suppose that

$$|k(t, s, x(s))| \leq f(t)g(s)|x(s)|, \quad (4)$$

$$|F(t, x(t), v)| \leq f(t)|x(t)| + |v|, \quad (5)$$

where  $f$  and  $g$  are real-valued nonnegative continuous functions defined on  $R^+$ . Then the solution  $x(t)$  of (1) is bounded and

$$|x(t)| \leq p_1(t) \left[ 1 + \int_0^t f(s) \exp \left( \int_0^s [f(\tau) + g(\tau)] d\tau \right) ds \right], \quad (6)$$

where

$$p_1(t) = |x_0| + \int_0^t |h(s)| ds,$$

for  $t \in R^+$ .

**Proof.** Multiplying both sides of equation (1) by  $x(t)$ , substituting  $t = s$  and integrating from 0 to  $t$  we have

$$x^2(t) = x_0^2 + 2 \int_0^t \left[ x(s) F \left( s, x(s), \int_0^s k(s, \tau, x(\tau)) d\tau \right) + h(s)x(s) \right] ds. \quad (7)$$

From (4),(5) we get

$$\begin{aligned} |x(t)|^2 &\leq |x_0|^2 + 2 \int_0^t \left[ |x(s)| \left( f(s)|x(s)| + \int_0^s |k(s, \tau, x(\tau))| d\tau \right) + |h(s)||u(s)| \right] ds \\ &\leq |x_0|^2 + 2 \int_0^t \left[ |x(s)| \left( f(s)|x(s)| + \int_0^s f(s)g(\tau)|x(\tau)| d\tau \right) + |h(s)||u(s)| \right] ds. \end{aligned}$$

Thus

$$|x(t)|^2 \leq |x_0|^2 + 2 \int_0^t \left[ f(s)|x(s)| \left( |x(s)| + \int_0^s g(\tau)|x(\tau)| d\tau \right) + |h(s)||x(s)| \right] ds. \quad (8)$$

The proof follows from Theorem 2.1 (I) and inequality (8).

### 3 Boundedness and asymptotic stability of solutions of integral equations

Consider the nonlinear Volterra integral equation of the form

$$u^p(t) = f(t) + \int_0^t k(t, s)g(s, u(s))ds, \quad (9)$$

where  $f : R_+ \rightarrow R$ ,  $k : R_+ \times R_+ \rightarrow R$ ,  $g : R_+ \times R \rightarrow R$  are continuous functions and  $p > 1$  is a constant. Okrasinski[5] studied the problem of existence and uniqueness of solutions of equation (9) in the form

$$u^p = k * u + f, \quad p > 1,$$

where  $k, f$  are known smooth functions depending on physical parameters. Pachpatte[6] investigated the boundedness and asymptotic behaviour of solutions of (9) using inequalities derived by himself in [6].

For an interesting discussion concerning the occurrence of equation (9) in the theory of water percolation phenomena and its physical meaning, see Okrasinski[5] and Pachpatte[10].

Now suppose  $u(t) \geq 0$ ,  $g \geq 0$ ,  $r_i \geq 0$ ,  $i = 1, 2, \dots, n-1$  are continuous functions defined on  $R_+$  and let  $p > 1$  be a constant.

Put

$$\begin{aligned} M[t, r, g(t_n)u(t_n)] &= M[t, r_1, \dots, r_{n-1}, g(t_n), u(t_n)] \\ &= \int_0^t r_1(t_1) \int_0^{t_1} r_2(t_2) \dots \int_0^{t_{n-2}} r_{n-1}(t_{n-1}) \int_0^{t_{n-1}} g(t_n)u(t_n) dt_n dt_{n-1} \dots dt_2 dt_1. \end{aligned}$$

In the following it is assumed that every solution  $u(t)$  of (9) exists on  $R_+$ .

For investigation of boundedness and asymptotic behaviour of solutions of (9) we use Pachpatte's iterated integral inequalities (see[10]) which we can formulate as follows:

**Theorem 3.1** . Assume

$$u^p(t) \leq c + M[t, r, g(t_n), u(t_n)],$$

for  $t \in R_+$  where  $c \geq 0$  is a constant.

Then

$$u(t) \leq [c^{(p-1)/p} + ((p-1)/p) M[t, r, g(t_n)]]^{1/(p-1)}.$$

By virtue of Theorem 3.1 we get the following results:

**Theorem 3.2** Suppose

$$|f(t)| \leq c_1, \quad |k(t, s)| \leq c_2, \quad |g(t, u)| \leq r(t)|u|, \quad (10)$$

where  $c_1, c_2$  are nonnegative constants,  $r : R_+ \rightarrow R_+$  is a continuous function. Then the solution  $u(t)$  of (9) is bounded and

$$|u(t)| \leq \left[ c_1^{(p-1)/p} + ((p-1)/p) \int_0^t c_2 r(s) ds \right]^{1/(p-1)}. \quad (11)$$

**Proof.** From (10) we obtain

$$|u(t)|^p \leq c_1 + \int_0^t c_2 r(s) |u(s)| ds$$

By Theorem 3.1 for  $n = 1$  we get inequality (11). The proof is complete.

**Theorem 3.3** Suppose

$$|f(t)| \leq ce^{-pt}, \quad |k(t, s)| \leq h(s)e^{-pt}, \quad |g(t, u)| \leq r(t)|u|, \quad (12)$$

where  $c$  is a nonnegative constant,  $r(t)$  is as defined above,  $h : R_+ \rightarrow R_+$  is a continuous function and

$$\int_0^\infty h(s)r(s)e^{-s} ds < \infty. \quad (13)$$

Then the solution  $u(t)$  of (9) is asymptotic stable.

**Proof.** From (12) we get

$$|u(t)|^p \leq c_1 e^{-pt} + \int_0^t h(s) e^{-pt} r(s) |u(s)| ds.$$

Then

$$(e^t |u(t)|)^p \leq c_1 + \int_0^t h(s) e^{-s} r(s) (e^s |u(s)|) ds$$

By Theorem 3.1 for  $n=1$  we obtain

$$e^t |u(t)| \leq \left[ (c_1)^{(p-1)/p} + (p-1)/p \int_0^t h(s) r(s) e^{-s} ds \right]^{1/(p-1)}.$$

Thus

$$|u(t)| \leq e^{-t} \left[ (c_1)^{(p-1)/p} + \frac{p-1}{p} K \right]^{1/(p-1)},$$

where  $K > 0$  is a constant which bounds integral (13).

Put

$$c^* = \left[ (c_1)^{(p-1)/p} + \frac{p-1}{p} K \right]^{1/(p-1)},$$

then

$$|u(t)| \leq c^* e^{-t}$$

which means that the solution  $u(t)$  of (9) approaches zero as  $t \rightarrow \infty$ . The proof is complete.

## Acknowledgement

This research has been supported by the Czech Ministry of Education in the frames of MSM002160503 Research Intention MIKROSYN New Trends in Microelectronic Systems and Nanotechnologies and MSM0021630529 Research Intention Intelligent Systems in Automatization.

## References

- [1] DRAGOMIR, S.S.: *On Volterra integral equations with kernels of L-type*, Ann. Univ. Timisoara Facult de Mat. Inform., vol.25, 21-41, 1987.
- [2] ENGLER, H.: *Global regular solutions for the dynamic antiplane shear problem in nonlinear viscoelasticity*, Math. Zentral. vol.202, 251-259, 1989.
- [3] FAJMON, B., ŠMARDÁ, Z.: *Application of integral inequalities*, Proceedings of the Seventh International Mathematical Workshop, FAST VUT Brno, 2008, 28-30.
- [4] HARAUX, A.: *Nonlinear Evolution Equations. Global Behavior of Solutions*, Lectures Notes in Math. No 847, Berlin, New York, Springer-Verlag, 1981.
- [5] OKRASINSKI, W.: *On a nonlinear convolution equation occurring in the theory of water precolation*, Ann. Polon. Math., vol.37, 223-229, 1980.

- [6] PACHPATTE, B.G.: *On a certain inequality arising in the theory of differential equations*, J. Math. Anal. Appl., vol.182, 143-157, 1994.
- [7] PACHPATTE, B.G.: *On certain nonlinear integral inequalities and their discrete analogues*, Facta Univ. (NIS), Ser. Math. Inform., vol. 8, 21-34, 1993.
- [8] PACHPATTE, B.G.: *On some fundamental integral inequalities arising in the theory of differential equations*, Chinese J. Math., vol. 22, 261-273, 1994.
- [9] PACHPATTE, B.G.: *On a new inequality suggested by the study of certain epidemic models*, J. Math. Anal. Appl., vol. 195, 638-644, 1995.
- [10] PACHPATTE, B.G.: *Inequalities for differential and integral equations*, Mathematics in Science and Engineering, vol.197, Academic Press Inc. , 2006.
- [11] YANG, E.H.: *On asymptotic behaviour of certain integro-differential equations*, Proc. Am. Math. Soc., vol. 90, 271- 276.

### **Current address**

#### **RNDr. Břetislav Fajmon, Ph.D.**

Department of Mathematics

Faculty of Electrical Engineering and Communication

Brno University of Technology, Technická 8, 616 00 BRNO

e-mail: fajmon@feec.vutbr.cz

#### **Doc. RNDr. Zdeněk Šmarda, CSc.**

Department of Mathematics

Faculty of Electrical Engineering and Communication

Brno University of Technology, Technická 8, 616 00 BRNO

e-mail: smarda@feec.vutbr.cz



# DISCRETE MAXIMUM PRINCIPLES FOR PARABOLIC PROBLEMS WITH GENERAL BOUNDARY CONDITIONS

FARAGÓ István, (H), HORVÁTH Róbert, (H), KOROTOV Sergey, (F)

**Abstract.** In this work, for the first time with respect to parabolic problems and discrete maximum principles, the cases with the mixed boundary conditions and an additional reactive term presented in the governing equation are considered. We derive the relevant continuous maximum principle, and also give its discrete analogue, when simplicial finite elements and the  $\theta$  time discretization method are used.

**Key words and phrases.** parabolic problem, maximum principle, linear finite elements, discrete maximum principle

*Mathematics Subject Classification.* 65M60, 65M50, 35B50

## 1 Introduction

Consider the following parabolic problem with general boundary conditions: Find a function  $u = u(t, x)$  such that

$$\frac{\partial u}{\partial t} - b\Delta u + cu = f \quad \text{in } Q_T := (0, T) \times \Omega, \quad (1)$$

$$u = g \quad \text{on } S_T^D := (0, T) \times \partial\Omega_D, \quad (2)$$

$$b\nabla u \cdot \nu = q \quad \text{on } S_T^N := (0, T) \times \partial\Omega_N, \quad (3)$$

$$b\nabla u \cdot \nu + \sigma u = r \quad \text{on } S_T^R := (0, T) \times \partial\Omega_R, \quad (4)$$

$$u|_{t=0} = u^0 \quad \text{in } \Omega, \quad (5)$$

where  $\Omega \subset \mathbf{R}^d$ ,  $d = 1, 2, 3, \dots$ , is a bounded polytopic domain with Lipschitz boundary  $\partial\Omega$ . Assume that  $\partial\Omega = \partial\Omega_D \cup \partial\Omega_N \cup \partial\Omega_R$ , where  $\partial\Omega_D \neq \emptyset$  and is closed,  $\partial\Omega_N$  and  $\partial\Omega_R$  are mutually disjoint measurable open sets. The subscripts (or superscripts)  $D$ ,  $N$ , and  $R$  always stand for Dirichlet, Neumann, and Robin types of boundary conditions, respectively,  $\nu$  is the outward normal to  $\partial\Omega$ ,  $T > 0$ , the problem coefficients are constant and such that

$$b > 0, \quad c \geq 0, \quad \sigma > 0, \quad q < 0, \quad (6)$$

and  $f, g, q, r, u^0$  are given functions. For any  $t \in (0, T)$  let  $Q_t$  stand for the cylinder  $(0, t) \times \Omega$ , and let  $\Gamma_0 := \{0\} \times \Omega$  denote its bottom. Moreover, let us define  $Q_{\bar{t}} := (0, t] \times \Omega$ ,  $S_{\bar{t}}^D := [0, t] \times \partial\Omega_D$ ,  $S_{\bar{t}}^N := [0, t] \times \partial\Omega_N$ , and  $S_{\bar{t}}^R := [0, t] \times \partial\Omega_R$ .

In the sequel we assume that all the given functions are sufficiently smooth so that the classical solution of problem (1)–(5) exists in the space  $C^{1,2}(Q_{\bar{T}}) \cap C^{0,1}(Q_{\bar{T}} \cup S_{\bar{T}}^D \cup S_{\bar{T}}^N \cup S_{\bar{T}}^R \cup \Gamma_0)$  and it is unique.

Then an upper bound for the solution can be given as follows [3]: For all  $t_1 \in (0, T)$ , it holds

$$u(t_1, x) \leq \max\{0; \max_{\Gamma_0 \cup S_{t_1}^D} u\} + \max\{0; \max_{S_{t_1}^R} \frac{r}{\sigma}\} + t_1 \max\{0; \max_{Q_{t_1}} f\}. \quad (7)$$

Inequality (7), under the assumptions (6), represents the form of the continuous maximum principle that we shall deal with for the above defined parabolic problem (1)–(5).

## 2 Linear finite element discretization

In the sequel we consider the following finite element (FE) discretization. Let

$$H_{\partial\Omega_D}^1(\Omega) = \{v \in H^1(\Omega) \mid v|_{\partial\Omega_D} = 0\}. \quad (8)$$

We assume that a simplicial partition  $\mathcal{T}_h$  of  $\bar{\Omega}$  is given, where  $h$  denotes the standard discretization parameter (the maximal diameter of elements from  $\mathcal{T}_h$ ), and that the partition is conforming and is such that any facet of any element is either a facet of the adjacent element or a part of the boundary. Let  $B_1, \dots, B_N$  denote all interior nodes and the nodes belonging to  $\partial\Omega_N \cup \partial\Omega_R$ , and let  $B_{N+1}, \dots, B_{\bar{N}}$  be the nodes lying on  $\partial\Omega_D$ . We also stand  $N_{\partial} := \bar{N} - N$ .

Let  $\phi_1, \dots, \phi_{\bar{N}}$  be the continuous piecewise linear nodal basis functions associated with nodes  $B_1, \dots, B_{\bar{N}}$ , respectively. It is obvious that

$$\phi_i \geq 0, \quad i = 1, \dots, \bar{N}, \quad \text{and} \quad \sum_{i=1}^{\bar{N}} \phi_i \equiv 1 \quad \text{in } \bar{\Omega}. \quad (9)$$

We denote the span of the basis functions by  $V^h \subset H^1(\Omega)$ , and define its subspace

$$V_{\partial\Omega_D}^h = \{v \in V^h \mid v|_{\partial\Omega_D} = 0\} \subset H_{\partial\Omega_D}^1(\Omega).$$

In what follows, we assume that the discretization of the initial and boundary conditions are linear interpolants in  $V^h$ , i.e.,

$$u_h^0(x) = \sum_{i=1}^{\bar{N}} u^0(B_i) \phi_i(x), \quad (10)$$

and

$$g_h(t, x) = \sum_{i=1}^{N_\partial} g_i^h(t) \phi_{N+i}(x), \quad \text{where} \quad g_i^h(t) = g(t, B_{N+i}), \quad i = 1, \dots, N_\partial. \quad (11)$$

From the consistency of the initial and the boundary conditions  $g(0, s) = u^0(s)$ ,  $s \in \partial\Omega_D$ , we have  $g_i^h(0) = u^0(B_{N+i})$ ,  $i = 1, \dots, N_\partial$ .

We search for a semidiscrete solution of the form

$$u_h(t, x) = \sum_{j=1}^N u_j^h(t) \phi_j(x) + g_h(t, x) = \sum_{j=1}^N u_j^h(t) \phi_j(x) + \sum_{j=N+1}^{\bar{N}} g_{j-N}^h(t) \phi_j(x). \quad (12)$$

Introducing the notation

$$\mathbf{v}^h(t) = [u_1^h(t), \dots, u_N^h(t), g_1^h(t), \dots, g_{N_\partial}^h(t)]^T, \quad (13)$$

we get a Cauchy problem for the systems of ordinary differential equations

$$\mathbf{M} \frac{d\mathbf{v}^h}{dt} + \mathbf{K} \mathbf{v}^h = \mathbf{f} + \mathbf{q} + \mathbf{r}, \quad \mathbf{v}^h(0) = [u^0(B_1), \dots, u^0(B_N), g_1^h(0), \dots, g_{N_\partial}^h(0)]^T \quad (14)$$

for the solution of the semidiscrete problem, where

$$\mathbf{M} = (m_{ij})_{i=1, j=1}^{N, \bar{N}}, \quad m_{ij} = \int_{\Omega} \phi_j \phi_i dx,$$

$$\mathbf{K} = (k_{ij})_{i=1, j=1}^{N, \bar{N}}, \quad k_{ij} = b \int_{\Omega} \nabla \phi_j \nabla \phi_i dx + c \int_{\Omega} \phi_j \phi_i dx + \sigma \int_{\partial\Omega_R} \phi_j \phi_i ds,$$

$$\mathbf{f} = [f_1, \dots, f_N]^T, \quad f_i = \int_{\Omega} f \phi_i dx,$$

$$\mathbf{q} = [q_1, \dots, q_N]^T, \quad q_i = \int_{\partial\Omega_N} q \phi_i ds,$$

and

$$\mathbf{r} = [r_1, \dots, r_N]^T, \quad r_i = \int_{\partial\Omega_R} r \phi_i ds.$$

In order to get a fully discrete numerical scheme, we choose a time-step  $\Delta t$  and denote the approximations to  $\mathbf{v}^h(n\Delta t)$ ,  $\mathbf{f}(n\Delta t)$ ,  $\mathbf{q}(n\Delta t)$ , and  $\mathbf{r}(n\Delta t)$  by  $\mathbf{v}^n$ ,  $\mathbf{f}^n$ ,  $\mathbf{q}^n$ , and  $\mathbf{r}^n$ ,  $n = 0, 1, \dots, n_T$  ( $n_T\Delta t = T$ ), respectively.

To discretize (14), we apply the  $\theta$ -method ( $\theta \in (0, 1]$  is a given parameter) and obtain a system of linear algebraic equations

$$\mathbf{M} \frac{\mathbf{v}^{n+1} - \mathbf{v}^n}{\Delta t} + \theta \mathbf{K} \mathbf{v}^{n+1} + (1 - \theta) \mathbf{K} \mathbf{v}^n = \mathbf{f}^{(n,\theta)} + \mathbf{q}^{(n,\theta)} + \mathbf{r}^{(n,\theta)}, \quad (15)$$

where  $\mathbf{f}^{(n,\theta)} := \theta \mathbf{f}^{n+1} + (1 - \theta) \mathbf{f}^n$ ,  $\mathbf{q}^{(n,\theta)} := \theta \mathbf{q}^{n+1} + (1 - \theta) \mathbf{q}^n$ , and  $\mathbf{r}^{(n,\theta)} := \theta \mathbf{r}^{n+1} + (1 - \theta) \mathbf{r}^n$ .

Further, (15) can be rewritten as

$$(\mathbf{M} + \theta \Delta t \mathbf{K}) \mathbf{v}^{n+1} = (\mathbf{M} - (1 - \theta) \Delta t \mathbf{K}) \mathbf{v}^n + \Delta t \mathbf{f}^{(n,\theta)} + \Delta t \mathbf{q}^{(n,\theta)} + \Delta t \mathbf{r}^{(n,\theta)}, \quad (16)$$

where  $n = 0, 1, \dots, n_T - 1$ , and  $\mathbf{v}^0 = \mathbf{v}^h(0)$ .

Let  $\mathbf{A} := \mathbf{M} + \theta \Delta t \mathbf{K}$  and  $\mathbf{B} := \mathbf{M} - (1 - \theta) \Delta t \mathbf{K}$ . We shall use the following partitions of the matrices and vectors:

$$\mathbf{A} = [\mathbf{A}_0 | \mathbf{A}_\partial], \quad \mathbf{B} = [\mathbf{B}_0 | \mathbf{B}_\partial], \quad \mathbf{v}^n = [(\mathbf{u}^n)^T | (\mathbf{g}^n)^T]^T, \quad (17)$$

where  $\mathbf{A}_0$  and  $\mathbf{B}_0$  are  $(N \times N)$  matrices,  $\mathbf{A}_\partial, \mathbf{B}_\partial$  are of size  $(N \times N_\partial)$ ,  $\mathbf{u}^n = [u_1^n, \dots, u_N^n]^T \in \mathbf{R}^N$  and  $\mathbf{g}^n = [g_1^n, \dots, g_{N_\partial}^n]^T \in \mathbf{R}^{N_\partial}$ . The iterative scheme (16) can now be rewritten as follows

$$\mathbf{A} \mathbf{v}^{n+1} = \mathbf{B} \mathbf{v}^n + \Delta t \mathbf{f}^{(n,\theta)} + \Delta t \mathbf{q}^{(n,\theta)} + \Delta t \mathbf{r}^{(n,\theta)}, \quad (18)$$

or

$$[\mathbf{A}_0 | \mathbf{A}_\partial] \begin{bmatrix} \mathbf{u}^{n+1} \\ \mathbf{g}^{n+1} \end{bmatrix} = [\mathbf{B}_0 | \mathbf{B}_\partial] \begin{bmatrix} \mathbf{u}^n \\ \mathbf{g}^n \end{bmatrix} + \Delta t \mathbf{f}^{(n,\theta)} + \Delta t \mathbf{q}^{(n,\theta)} + \Delta t \mathbf{r}^{(n,\theta)}. \quad (19)$$

### 3 The discrete maximum principle

Let us define the following values for  $n = 0, \dots, n_T$ :

$$g_{max}^n = \max\{0, g_1^n, \dots, g_{N_\partial}^n\}, \quad (20)$$

$$v_{max}^n = \max\{0, g_{max}^n, u_1^n, \dots, u_N^n\}, \quad (21)$$

$$f_{max}^{(n,n+1)} = \max\{0, \max_{x \in \Omega, \tau \in (n\Delta t, (n+1)\Delta t)} f(\tau, x)\}, \quad (22)$$

$$r_{max}^{(n,n+1)} = \max\{0, \max_{x \in \partial\Omega_R, \tau \in (n\Delta t, (n+1)\Delta t)} r(\tau, x)\}, \quad (23)$$

for  $n = 0, \dots, n_T - 1$ .

Then the discrete maximum principle (DMP) corresponding to (7) (under condition  $q < 0$ ) takes the following form (cf. [4, p. 100]):

$$u_i^{n+1} \leq \max\{0, g_{max}^{n+1}, v_{max}^n\} + \frac{1}{\theta\sigma} r_{max}^{(n,n+1)} + \Delta t f_{max}^{(n,n+1)}, \quad (24)$$

for  $i = 1, \dots, N$ ;  $n = 0, \dots, n_T - 1$ .

We can give an algebraic condition for DMP as follows [2]:

**Theorem 3.1** *Galerkin approximation for the solution of problem (1)–(5), combined with the  $\theta$ -method for time discretization (where  $\theta \in (0, 1]$ ), satisfies the discrete maximum principle (24) under condition  $q < 0$  if*

$$\mathbf{A}_0^{-1} \geq \mathbf{0}, \quad (C1)$$

$$\mathbf{A}_0^{-1} \mathbf{A}_\partial \leq \mathbf{0}, \quad (C2)$$

$$\mathbf{A}_0^{-1} \mathbf{B} \geq \mathbf{0}. \quad (C3)$$

**Remark 3.2** *Conditions (C1)–(C3) are ensured by the following simpler assumptions*

$$\mathbf{A}_0^{-1} \geq \mathbf{0}, \quad (C1^*)$$

$$\mathbf{A}_\partial \leq \mathbf{0}, \quad (C2^*)$$

$$\mathbf{B} \geq \mathbf{0}. \quad (C3^*)$$

**Theorem 3.3** *Galerkin approximation for the solution of problem (1)–(5), combined with the  $\theta$ -method for time discretization (where  $\theta \in (0, 1]$ ), satisfies the discrete maximum principle (24) if*

$$k_{ij} \leq 0, \quad i = 1, \dots, N, \quad j = 1, \dots, \bar{N}, \quad i \neq j, \quad (C1')$$

$$m_{ij} + \theta \Delta t k_{ij} \leq 0, \quad i = 1, \dots, N, \quad j = 1, \dots, \bar{N}, \quad i \neq j, \quad (C2')$$

$$m_{ii} - (1 - \theta) \Delta t k_{ii} \geq 0, \quad i = 1, \dots, N. \quad (C3')$$

For the proofs of the above two theorems see, e.g. [4].

#### 4 DMP on simplicial meshes

We shall generally denote any simplex from  $\mathcal{T}_h$  by the symbol  $K$  and also use denotation  $\alpha_{ij}^K$  for the angle between  $(d-1)$ -dimensional facets  $F_i^K$  and  $F_j^K$  of  $K$  which is opposite to the edge connecting vertices  $B_i$  and  $B_j$ , and let  $h_i^K(h_j^K)$  be the height of  $K$  from  $B_i(B_j)$  onto  $F_i^K(F_j^K)$ .

The contributions to the mass matrix  $\mathbf{M}$  over the simplex  $K$  are (cf. [1])

$$m_{ij}|_K = \int_K \phi_i \phi_j dx = (1 + \delta_{ij}) \frac{d!}{(d+2)!} \text{meas}_d K, \quad (25)$$

where  $\delta_{ij}$  is Kronecker's symbol.

In order to compute the entries of the matrix  $\mathbf{K}$ , we shall use the following formulae presented e.g. in [1]:

$$\nabla \phi_i \cdot \nabla \phi_j|_K = -\frac{\text{meas}_{d-1} F_i^K \cdot \text{meas}_{d-1} F_j^K}{(d \text{meas}_d K)^2} \cos \alpha_{ij}^K = -\frac{\cos \alpha_{ij}^K}{h_i^K h_j^K} \quad (i \neq j), \quad (26)$$

$$\nabla \phi_i \cdot \nabla \phi_i|_K = \frac{(\text{meas}_{d-1} F_i^K)^2}{(d \text{meas}_d K)^2} = \frac{1}{(h_i^K)^2}. \quad (27)$$

Then the conditions  $(C1')$ – $(C3')$  in the Theorem 3.3 can be guaranteed by the following three lemmas.

**Lemma 4.1** *Let the simplicial partition  $\mathcal{T}_h$  of  $\bar{\Omega}$  be such that for any pair of distinct  $(d-1)$ -dimensional facets  $F_i^K$  and  $F_j^K$  of any simplex  $K$  from  $\mathcal{T}_h$ , we have*

$$\frac{\cos \alpha_{ij}^K}{h_i^K h_j^K} \geq \frac{c}{b(d+1)(d+2)} + \frac{\sigma}{bd(d+1)} \frac{\text{meas}_{d-1}(\partial K \cap \partial \Omega_R)}{\text{meas}_d K}, \quad (28)$$

where  $\partial K$  is the boundary of  $K$ . Then

$$k_{ij} \leq 0, \text{ for } i = 1, \dots, N, j = 1, \dots, \bar{N}, i \neq j.$$

**Lemma 4.2** *Let the simplicial partition  $\mathcal{T}_h$  of  $\bar{\Omega}$  and the time-step  $\Delta t$  be such that for any pair of distinct  $(d-1)$ -dimensional facets  $F_i^K$  and  $F_j^K$  of any simplex  $K$  from  $\mathcal{T}_h$ , we have*

$$\frac{\cos \alpha_{ij}^K}{h_i^K h_j^K} \geq \frac{c + 1/(\theta \Delta t)}{b(d+1)(d+2)} + \frac{\sigma}{bd(d+1)} \frac{\text{meas}_{d-1}(\partial K \cap \partial \Omega_R)}{\text{meas}_d K}. \quad (29)$$

where  $\partial K$  is the boundary of  $K$ . Then

$$m_{ij} + \theta \Delta t k_{ij} \leq 0, \text{ for } i = 1, \dots, N, j = 1, \dots, \bar{N}, i \neq j.$$

**Lemma 4.3** *Let the simplicial partition  $\mathcal{T}_h$  of  $\bar{\Omega}$  and the time-step  $\Delta t$  be such that for any simplex  $K$  from  $\mathcal{T}_h$ , we have*

$$0 \leq -\frac{1}{(h_i^K)^2} + \frac{2\left(\frac{1}{(1-\theta)\Delta t} - c\right)}{b(d+1)(d+2)} - \frac{2\sigma}{bd(d+1)} \frac{\text{meas}_{d-1}(\partial K \cap \partial\Omega_R)}{\text{meas}_d K}, \quad (30)$$

where  $\partial K$  is the boundary of  $K$ . Then

$$m_{ii} - (1 - \theta)\Delta t k_{ii} \geq 0, \quad i = 1, \dots, N.$$

Summarizing the above results we can formulate the main result of the paper.

**Theorem 4.4** *Galerkin approximation for the solution of problem (1)–(5), combined with the  $\theta$ -method for time discretization (where  $\theta \in (0, 1]$ ), satisfies the discrete maximum principle (24) if an acute simplicial mesh is used and the time-step satisfies the following (lower and upper) estimates:*

$$\frac{1}{\Delta t} \leq \theta \left( \frac{\cos \alpha_{ij}^K b(d+1)(d+2)}{h_i^K h_j^K} - \frac{\sigma(d+2)\text{meas}_{d-1}(\partial K \cap \partial\Omega_R)}{d \text{meas}_d K} - c \right), \quad (31)$$

and

$$\Delta t \leq \frac{1}{(1 - \theta) \left( \frac{b(d+1)(d+2)}{2(h_i^K)^2} + \frac{\sigma(d+2)\text{meas}_{d-1}(\partial K \cap \partial\Omega_R)}{d \text{meas}_d K} + c \right)}. \quad (32)$$

## Acknowledgement

István Faragó was supported by the Grant no. K67819 and T049819 of OTKA.

Róbert Horváth was supported by the Grant no. K67819 and K61800 of OTKA, and János Bolyai scholarship.

Sergey Korotov was supported by grant no. 127031 from the Academy of Finland.

## References

- [1] BRANDTS, J., KOROTOV, S., KŘÍŽEK, M., *Dissection of the Path-Simplex in  $\mathbf{R}^n$  into  $n$  Path-Subsimplices*, Linear Algebra Appl. 421 (2007), pp. 382–393.
- [2] FARAGÓ, I., HORVÁTH, R., KOROTOV, S., *Discrete Maximum Principle for Linear Parabolic Problems Solved on Hybrid Meshes*, Appl. Numer. Math. 53 (2005), pp. 249–264.
- [3] FARAGÓ, I., HORVÁTH, R., KOROTOV, S., *Discrete Maximum Principles for FE Solutions of Nonstationary Diffusion-Reaction Problems with Mixed Boundary Conditions*, Preprint A550 (2008), Helsinki University of Technology (submitted).

- [4] FUJII, H., *Some Remarks on Finite Element Analysis of Time-Dependent Field Problems*, Theory and Practice in Finite Element Structural Analysis, Univ. Tokyo Press, Tokyo (1973), pp. 91–106.

**Current address**

**István Faragó**

Department of Applied Analysis, Eötvös Loránd University  
H-1518, Budapest, Pf. 120, Hungary,  
e-mail: faragois@cs.elte.hu

**Róbert Horváth**

Institute of Mathematics and Statistics, University of West-Hungary  
Erzsébet u. 9, H-9400, Sopron, Hungary,  
e-mail: rhorvath@ktk.nyme.hu

**Sergey Korotov**

Institute of Mathematics, Helsinki University of Technology  
P.O. Box 1100, FIN-02015 TKK, Finland  
e-mail: sergey.korotov@hut.fi



## CHAOTIC OSCILLATIONS OF ELASTIC BEAMS

FEČKAN Michal, (SK)

**Abstract.** We survey our recent results on the existence of chaotic oscillations of weakly damped and periodically forced elastic beams. The bifurcation theory of chaotic oscillations is developed with several applications to concrete beam partial differential equations.

**Key words and phrases.** Differential equations, Homoclinic solutions, Bifurcations, Chaos, Elastic Beams.

*Mathematics Subject Classification.* Primary 34C37, 35B99; Secondary 74H65.

## 1 Introduction

A model for oscillations of an elastic beam with a compressive axial load  $P_0$  (see Figure 1) is given by the partial differential equation (PDE)

$$\ddot{u} = -u'''' - P_0 u'' + \left[ \int_0^\pi u'(s, t)^2 ds \right] u'' - 2\mu_2 \dot{u} + \mu_1 \cos \omega_0 t \quad (1)$$

where  $P_0$ ,  $\mu_1$ ,  $\mu_2$ ,  $\omega_0$  are constants and  $u$  is a real valued function of two variables  $t \in \mathbb{R}$ ,  $x \in [0, \pi]$ , subject to the boundary conditions

$$u(0, t) = u(\pi, t) = u''(0, t) = u''(\pi, t) = 0.$$

In (1), a superior dot denotes differentiation with respect to  $t$  and prime differentiation with respect to  $x$ . When  $P_0$  is sufficiently large, (1) can exhibit chaotic behavior.

In (1) substitute  $u(x, t) = \sum_{k=1}^{\infty} u_k(t) \sin kx$ , multiply by  $\sin nx$  and integrate from 0 to  $\pi$ .

This yields the infinite set of ordinary differential equations (ODEs)

$$\ddot{u}_n = n^2(P_0 - n^2)u_n - \frac{\pi}{2}n^2 \left[ \sum_{k=1}^{\infty} k^2 u_k^2 \right] u_n - 2\mu_2 \dot{u}_n + 2\mu_1 \left[ \frac{1 - (-1)^n}{\pi n} \right] \cos \omega_0 t,$$

$$n = 1, 2, \dots$$

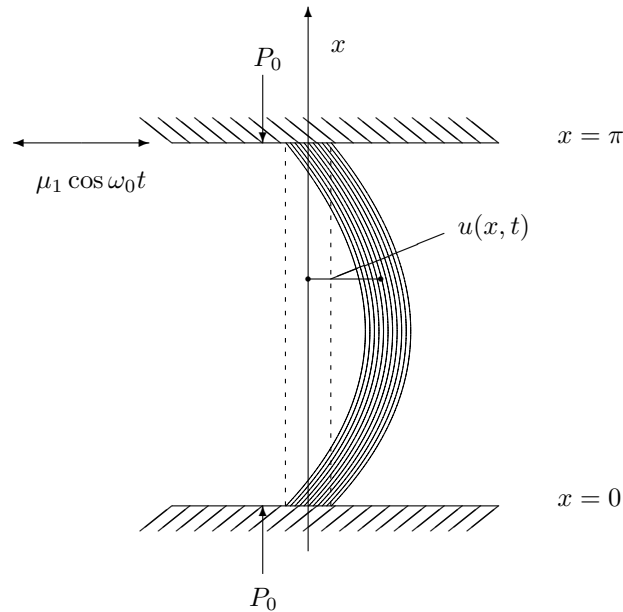


Figure 1: The forced buckled beam (1).

We see that the linear parts of these equations are uncoupled and the equations divide into two types. The system of equations defined by  $1 \leq n^2 < P_0$  has a hyperbolic equilibrium at the origin whereas, for the system of equations satisfying  $n^2 \geq P_0$ , this equilibrium is a center. For simplicity let us assume  $1 < P_0 < 4$ . Then only the equation with  $n = 1$  is hyperbolic while the system of remaining equations has a center. To emphasize this let us put  $x = (u_1, \dot{u}_1)$  and

$$y = (u_2, \dot{u}_2/\omega_1, u_3, \dot{u}_3/\omega_2, \dots),$$

where we have defined  $a^2 = P_0 - 1$  and  $\omega_n^2 = (n+1)^2 [(n+1)^2 - P_0]$ . The preceding equations now take the form

$$\begin{aligned} \dot{x}_1 &= x_2, \\ \dot{x}_2 &= a^2 x_1 - \frac{\pi}{2} \left[ x_1^2 + \sum_{k=1}^{\infty} (k+1)^2 y_{2k-1}^2 \right] x_1 \\ &\quad - 2\mu_2 x_2 + \frac{4}{\pi} \mu_1 \cos \omega_0 t, \\ \dot{y}_{2n-1} &= \omega_n y_{2n}, \\ \dot{y}_{2n} &= -\omega_n y_{2n-1} - \frac{\pi (n+1)^2}{2 \omega_n} \left[ x_1^2 + \sum_{k=1}^{\infty} (k+1)^2 y_{2k-1}^2 \right] y_{2n-1} \\ &\quad - 2\mu_2 y_{2n} + 2\mu_1 \left[ \frac{1 - (-1)^{n+1}}{\pi (n+1) \omega_n} \right] \cos \omega_0 t. \end{aligned} \tag{2}$$

In (2) we project onto the hyperbolic subspace by setting  $y = 0$  to obtain what we shall call the reduced equation. In our example this is

$$\ddot{x}_1 = a^2 x_1 - \frac{\pi}{2} x_1^3 - 2\mu_2 \dot{x}_1 + \frac{4}{\pi} \mu_1 \cos \omega_0 t. \tag{3}$$

We see that this is the forced, damped Duffing equation with negative stiffness for which standard theory yields chaotic dynamics. The purpose of the present work is to survey results of [7] where it is shown that the chaotic dynamics of (3) is, in some sense, shadowed in the dynamics of the full equation (2).

It is interesting to look at some history of this problem. The first work was by Holmes [9] in which he started with the PDE and carried out the Galerkin expansion but restricted his analysis to the reduced equation (3). The significance of that work is that it introduced the idea of Melnikov analysis. In subsequent work [10] Holmes and Marsden extended the results to infinite dimension but abandoned the Galerkin approach in favor of nonlinear semigroup techniques directly in infinite dimensions. In our work we go back to the original, simpler analysis of the reduced equation and then show that the results apply to the original PDE. Some advantages to this are that the Galerkin projection is a technique familiar to many engineers and physicists and, also, we are able to utilize our general Melnikov results. This is illustrated further in the generalizations which follow.

Next, we show the chaos in [1] for elastic beams of the form

$$\begin{aligned}\ddot{u} + u'''' + \varepsilon \delta \dot{u} + \varepsilon \mu h(x, \sqrt{\varepsilon} t) &= 0, \\ u''(0, \cdot) &= u''(\pi/4, \cdot) = 0, \\ u'''(0, \cdot) &= -\varepsilon f(u(0, \cdot)), \quad u'''(\pi/4, \cdot) = \varepsilon f(u(\pi/4, \cdot)),\end{aligned}\tag{4}$$

where  $\varepsilon > 0$  and  $\mu$  are sufficiently small parameters,  $\delta > 0$  is a constant,  $f \in C^2(\mathbb{R})$ ,  $h \in C^2([0, \pi/4] \times \mathbb{R})$  and  $h(x, t)$  is 1-periodic in  $t$ , provided an associated reduced equation has a homoclinic orbit. Equation (4) describes vibrations of a beam resting on two identical bearings with purely elastic responses which are determined by  $f$ . The length of the beam is  $\pi/4$ . Since  $\varepsilon > 0$ , (4) is a semilinear problem.

On the other hand, we study the existence of chaos in [2] for the following PDE

$$\begin{aligned}\ddot{u} + u'''' + (i^2 + \varepsilon \sigma^2)u'' - \varepsilon \kappa u'' f\left(\int_0^\pi u'(s, t)^2 ds\right) &= \varepsilon(\nu h(x, \sqrt{\varepsilon} t) - \delta \dot{u}), \\ u(0, t) &= u(\pi, t) = 0 = u''(0, t) = u''(\pi, t),\end{aligned}$$

where  $\kappa > 0$ ,  $\delta > 0$  and  $\sigma \in (0, 1]$  are constants,  $\varepsilon > 0$  and  $\nu$  are small parameters,  $i \in \mathbb{N}$  is fixed,  $h(x, t)$  is periodic in time. Here the external load  $i^2 + \varepsilon \sigma^2$  is resonant and the contribution given from the stress due to the external rigidity  $\varepsilon \kappa$ , does not drive the system too far away from the resonance.

Finally, some more recent work on the chaos in PDEs is by Berti and Carminati [4]. An undamped buckled beam is investigated by Yagasaki [15] to show Arnold diffusion type motions. Perturbed nonlinear Schrödinger equations are studied by Li [11, 12] under generic conditions.

## 2 The Abstract Problem

Using the example in the preceding section as a model we now develop an abstract theory. Let  $\mathbb{Y}$  and  $\mathbb{H}$  be separable real Hilbert spaces with  $\mathbb{Y} \subset \mathbb{H}$ .

We now consider differential equations of the form

$$\begin{aligned}\dot{x} &= f(x, y, \mu, t) = f_0(x, y) + \mu_1 f_1(x, y, \mu, t) + \mu_2 f_2(x, y, \mu, t), \\ \dot{y} &= g(x, y, \mu, t) = Ay + g_0(x, y) + \mu_1 \nu \cos \omega_0 t + \mu_2 g_2(x, y, \mu),\end{aligned}\tag{5}$$

with  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{Y}$ ,  $\mu = (\mu_1, \mu_2) \in \mathbb{R}^2$ ,  $\nu \in \mathbb{Y}$ . We make the following assumptions about (5):

(H1)  $A : \mathbb{Y} \rightarrow \mathbb{H}$  is a continuous and linear transformation.

(H2) The functions  $f_i$  and  $g_i$  are in the spaces:

$$\begin{aligned}f_0 &\in \mathcal{C}^4(\mathbb{R}^n \times \mathbb{Y}, \mathbb{R}^n); & f_1, f_2 &\in \mathcal{C}^4(\mathbb{R}^n \times \mathbb{Y} \times \mathbb{R}^2 \times \mathbb{R}, \mathbb{R}^n); \\ g_0 &\in \mathcal{C}^4(\mathbb{R}^n \times \mathbb{Y}, \mathbb{Y}); & g_2 &\in \mathcal{C}^4(\mathbb{R}^n \times \mathbb{Y} \times \mathbb{R}^2, \mathbb{Y}).\end{aligned}$$

(H3)  $f_1$  and  $f_2$  are periodic in  $t$  with period  $T = 2\pi/\omega_0$ .

(H4)  $f_0(0, 0) = 0$  and  $D_2 f_0(x, 0) = 0$ .

(H5) The eigenvalues of  $D_1 f_0(0, 0)$  lie off the imaginary axis.

(H6) The equation  $\dot{x} = f_0(x, 0)$  has a nontrivial solution homoclinic to  $x = 0$ .

(H7)  $g_0(x, 0) = g_2(x, 0, \mu) = 0$ ,  $D_{12} g_0(0, 0) = 0$  and  $D_{22} g_0(x, 0) = 0$ .

(H8) There are constants  $K > 0$ ,  $\delta > 0$  and  $b > 0$  so that when  $0 \leq |\mu_2| \leq \delta$  the variational equation

$$\dot{v} = (A + \mu_2 D_2 g_2(0, 0, 0))v$$

has a group  $\{V_{\mu_2}(t)\}$  of bounded evolution operators from  $\mathbb{Y}$  to  $\mathbb{Y}$  satisfying  $|V_{\mu_2}(t)V_{\mu_2}(s)^{-1}| \leq K e^{b\mu_2(s-t)}$ .

(H9) There is a constant  $K > 0$  such that the nonhomogeneous variational equation

$$\dot{v} = [A + \mu_2 D_2 g_2(0, 0, 0)]v + \mu_1 \nu \cos \omega_0 t$$

has a particular solution  $\psi : \mathbb{R} \rightarrow \mathbb{Y}$  satisfying  $|\psi(t)| \leq K|\mu_1||\nu|$ .

By a weak solution to (5) we mean a pair of continuous functions  $x_0 : \mathbb{R} \rightarrow \mathbb{R}^n$ ,  $y_0 : \mathbb{R} \rightarrow \mathbb{Y}$  such that  $x_0$  is differentiable and  $y_0$  has a derivative  $\dot{y}_0 : \mathbb{R} \rightarrow \mathbb{H}$  and which satisfy (5) pointwise in  $\mathbb{H}$ .

By (H8) we mean that  $V_{\mu_2}(s)^{-1} = V_{\mu_2}(-s)$ ,  $V_{\mu_2}(s) \circ V_{\mu_2}(t) = V_{\mu_2}(s+t)$ ,  $V_{\mu_2}(0) = \mathbb{I}$  and that for  $y_0 \in \mathbb{Y}$ ,  $y(t) = V_{\mu_2}(t)y_0$  is the weak solution to  $\dot{v} = [A + \mu_2 D_2 g_2(0, 0, 0)]v$  satisfying  $y(0) = y_0$ .

### 3 Chaotic Oscillations of the Abstract Problem

The reduced system of equations for (5) is

$$\dot{x} = f(x, 0, \mu, t) = f_0(x, 0) + \mu_1 f_1(x, 0, \mu, t) + \mu_2 f_2(x, 0, \mu, t) \quad (6)$$

with  $x \in \mathbb{R}^n$ . In [3, 5, 6, 8, 13] a general Melnikov theory is developed for first order systems in  $\mathbb{R}^n$ . We summarize those results here as applied to (6).

By (H6), (6) has a nontrivial homoclinic solution  $\gamma$  when  $\mu = 0$ . By the variational equation along  $\gamma$  we mean the linear equation

$$\dot{u} = D_1 f_0(\gamma, 0)u \quad (7)$$

and by the adjoint the system

$$\dot{v} = -D_1 f_0(\gamma, 0)^* v. \quad (8)$$

We let  $\{u_1, \dots, u_d\}$  denote a basis for the vector space of bounded solutions to (7) with  $u_d = \dot{\gamma}$  and we let  $\{v_1, \dots, v_d\}$  denote a basis for the vector space of bounded solutions to (8). Now define the functions  $a_{ij} : \mathbb{R} \rightarrow \mathbb{R}$ , constants  $b_{ijk}$  and function

$$M : \mathbb{R}^2 \times \mathbb{R} \times \mathbb{R}^{d-1} \rightarrow \mathbb{R}^d$$

by

$$\begin{aligned} a_{ij}(\alpha) &= \int_{-\infty}^{\infty} \langle v_i(t), f_j(\gamma(t), 0, 0, t + \alpha) \rangle dt; \\ i &= 1, \dots, d; \quad j = 1, 2; \\ b_{ijk} &= \int_{-\infty}^{\infty} \langle v_i, D_{11} f_0(\gamma, 0) u_j u_k \rangle dt; \\ i &= 1, \dots, d; \quad j, k = 1, \dots, d-1; \end{aligned} \quad (9)$$

$$M_i(\mu, \alpha, \beta) = \sum_{j=1}^2 a_{ij}(\alpha) \mu_j + \frac{1}{2} \sum_{j,k=1}^{d-1} b_{ijk} \beta_j \beta_k; \quad 1 \leq i \leq d.$$

The function  $M$  is our bifurcation function.

Now, suppose that (6) has a  $(d-1)$ -parameter family of homoclinic orbits given by  $t \rightarrow \gamma_\beta(t)$  with  $\beta \in U_0$  where  $U_0$  is an open neighborhood of the origin in  $\mathbb{R}^{d-1}$ . Then in (9) all  $b_{ijk} = 0$  and an alternate bifurcation function is required.

For each fixed  $\beta$  we let  $\{v_{\beta 1}, \dots, v_{\beta d}\}$  denote a basis for the vector space of bounded solutions to the adjoint equation  $\dot{v} = -D_1 f_0(\gamma_\beta, 0)^* v$ . Without loss of generality we can assume that each  $v_{\beta i}$  depends differentially on  $\beta$ . Now define functions  $a_{ij} : \mathbb{R} \times U_0 \rightarrow \mathbb{R}$  and  $M : \mathbb{R}^2 \times \mathbb{R} \times U_0 \rightarrow \mathbb{R}^d$  by

$$\begin{aligned} a_{ij}(\alpha, \beta) &= \int_{-\infty}^{\infty} \langle v_{\beta i}(t), f_j(\gamma_\beta(t), 0, 0, t + \alpha) \rangle dt; \\ i &= 1, \dots, d; \quad j = 1, 2; \end{aligned} \quad (10)$$

$$M_i(\mu, \alpha, \beta) = \sum_{j=1}^2 a_{ij}(\alpha, \beta) \mu_j; \quad 1 \leq i \leq d.$$

This function,  $M$ , is the bifurcation function for this situation. Now we can state the main theorem (cf. [7]).

**Theorem 3.1** *Suppose (H1)-(H10) hold. Let  $M$  be as in (9) or (10) and suppose  $(\mu_0, \alpha_0, \beta_0)$  are such that  $M(\mu_0, \alpha_0, \beta_0) = 0$  and  $D_{(\alpha, \beta)}M(\mu_0, \alpha_0, \beta_0)$  is nonsingular. Then  $\exists \bar{\xi}_0 > 0$  such that for any  $\mu = \xi\mu_0$ ,  $0 < \xi \leq \bar{\xi}_0$ , (5) possesses a countable infinity of subharmonic solutions of all possible periods, an uncountable infinity of chaotic solutions and it has a sensitive dependence on initial conditions.*

In order to apply the above results, we use for a concrete PDE the following procedure:

1. Use a Galerkin expansion to convert the PDE to an infinite set of ODEs as (5).
2. Truncate the equation to get the finite problem (6) and derive the Melnikov functions either (9) or (10). For this we must verify (H1) through (H6).
3. Use Theorem 3.1 to show chaos for the original PDE. This requires (H7)-(H9).

**Remark 3.2** *We know from [7] that in Theorem 3.1, (5) has a Smale horseshoe with the corresponding deterministic chaos (cf. [13, 14]).*

## 4 Applications to Elastic Beams

We apply the above procedure to a number of different cases and generalizations of the example in Section 1.

### 4.1 Planer Motion with One Buckled Mode

The boundary value problem for planer deflections of an elastic beam with a compressive axial load  $P_0$  and pinned ends is given by (1). The Melnikov function (9) with  $d = 1$  (cf. [7]) becomes

$$M(\alpha) = \left[ \frac{8\omega_0}{\sqrt{\pi}} \sin \omega_0 \alpha \operatorname{sech} \frac{\pi\omega_0}{2a} \right] \mu_1 - \left( \frac{16a^3}{3\pi} \right) \mu_2.$$

Thus, we obtain the following result using Theorem 3.1.

**Theorem 4.1** *Suppose  $1 < P_0 < 4$ . If  $\omega_0 \neq \omega_n \forall n \in \mathbb{N}$  then for any  $\mu_1$  and  $\mu_2 \neq 0$  small satisfying*

$$|\mu_2| < \frac{3\sqrt{\pi}\omega_0}{2a^3} \operatorname{sech} \frac{\pi\omega_0}{2a} \cdot |\mu_1|, \quad (11)$$

*PDE (1) has a homoclinic solution with the associated chaos.*

In the  $\mu_1$ - $\mu_2$  plane we get from the condition (11) four small open wedge-shaped regions of parameter values for which (1) exhibits chaos (see Figure 2). These regions are bounded by the lines  $\mu_1/\mu_2 = \pm \frac{3\sqrt{\pi}\omega_0}{2a^3} \operatorname{sech} \frac{\pi\omega_0}{2a}$  and  $\mu_2 = 0$ .

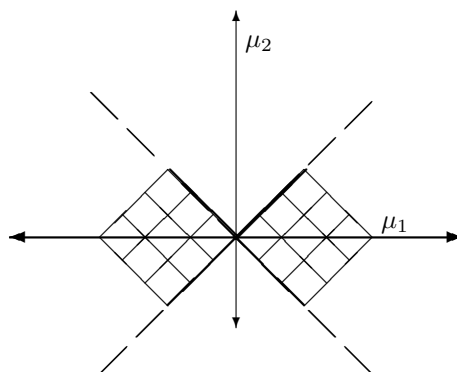


Figure 2: The chaotic open wedge-shaped region of (1) in  $\mathbb{R}^2$ .

## 4.2 Nonplanar Motion of a Symmetric Beam with One Buckled Mode

Let us consider a beam with symmetric cross section, pinned ends and compressive axial load  $P_0$  and assume now that the beam is not constrained to deflect in a plane. If  $u(x, t)$  and  $w(x, t)$  denote the transverse deflections at position  $x$  and time  $t$  we obtain the following boundary value problem.

$$\begin{aligned} \ddot{u} &= -u'''' - P_0 u'' + \left[ \int_0^\pi (u'(s, t)^2 + w'(s, t)^2) ds \right] u'' \\ &\quad - 2\mu_2 \dot{u} \cos \eta + \mu_1 \cos \zeta \cos \omega_0 t, \\ \ddot{w} &= -w'''' - P_0 w'' + \left[ \int_0^\pi (u'(s, t)^2 + w'(s, t)^2) ds \right] w'' \\ &\quad - 2\mu_2 \dot{w} \sin \eta + \mu_1 \sin \zeta \cos \omega_0 t, \\ u(0, t) &= u(\pi, t) = u''(0, t) = u''(\pi, t) = w(0, t) \\ &= w(\pi, t) = w''(0, t) = w''(\pi, t) = 0 \end{aligned} \quad (12)$$

where  $\eta, \zeta$  are constants. The parameters  $\mu_1, \mu_2$  represent the coefficients of, respectively, total transverse forcing and total viscous damping. These effects are distributed between the two directions of motion. The quantity  $\tan \zeta$  represents the ratio of forcing in the  $u$ -direction to forcing in the  $w$ -direction while  $\tan \eta$  plays the same role for the damping. We suppose  $\eta, \zeta \in (0, \pi/2)$  in order to avoid certain degeneracies and  $1 < P_0 < 4$ .

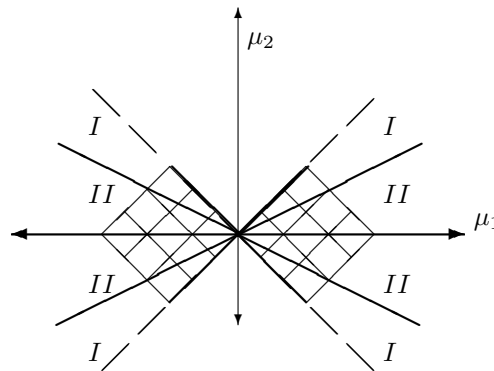


Figure 3: The chaotic wedge-shaped regions of (12) in  $\mathbb{R}^2$ .

Now the Melnikov function (10) has the form

$$M_1(\mu, \alpha, \beta) = \left[ \frac{8}{\sqrt{\pi}} \sin(\beta - \zeta) \cos \omega_0 \alpha \operatorname{sech} \frac{\pi \omega_0}{2a} \right] \mu_1,$$

$$M_2(\mu, \alpha, \beta) = \left[ \frac{8\omega_0}{\sqrt{\pi}} \cos(\beta - \zeta) \sin \omega_0 \alpha \operatorname{sech} \frac{\pi \omega_0}{2a} \right] \mu_1$$

$$- \left[ \frac{16a^3 (\cos \eta \cos^2 \beta + \sin \eta \sin^2 \beta)}{3\pi} \right] \mu_2.$$

The following result is obtained from Theorem 3.1 (cf. [7]).

**Theorem 4.2** Suppose  $\omega_0 \neq \omega_n$  for all  $n$  and let

$$m_1 = \frac{3\sqrt{\pi}\omega_0}{2a^2 (\cos \eta \cos^2 \zeta + \sin \eta \sin^2 \zeta)} \operatorname{sech} \frac{\pi \omega_0}{2a},$$

$$m_2 = \max_{\beta \in \mathbb{R}} \left\{ \frac{3\omega_0 \sqrt{\pi}}{2a^3} \frac{\cos(\beta - \zeta)}{\cos \eta \cos^2 \beta + \sin \eta \sin^2 \beta} \operatorname{sech} \frac{\pi \omega_0}{2a} \right\}.$$

- i) If  $m_0 \neq 0$  satisfies one but not both of  $|m_0| < m_i$  then if  $\mu_2 = m_0 \mu_1$  for  $\mu_1 \neq 0$  and  $\mu_2$  small then there exist two homoclinic orbits of (12) with the associated chaos.
- ii) If  $m_0 \neq 0$  satisfies each of  $|m_0| < m_i$  then there are four homoclinic orbits of (12) with the associated chaos.

Summarizing, we obtain eight open small wedge-shaped regions of parameter values in the  $\mu_1$ - $\mu_2$  plane bounded by the lines  $\mu_2/\mu_1 = \pm m_1$ ,  $\mu_2/\mu_1 = \pm m_2$  and  $\mu_2 = 0$  with  $m_1 \leq m_2$  for which the partial differential equation exhibits chaos (see Figure 3). In the regions labeled I there are two homoclinics while in regions II there exist four. It is interesting to note that in this case, by adjusting the parameters  $\eta$  and  $\zeta$ , it is possible to make the size of the wedge arbitrarily close to filling the  $\mu_1$ - $\mu_2$  plane.



### 4.3 Nonplanar, Nonsymmetric Beam with One Buckled Mode in Each Plane

For the case of a nonsymmetric beam with nonplanar motion we have the boundary value problem

$$\begin{aligned} \ddot{u} &= -u'''' - P_0 u'' + \left[ \int_0^\pi (u'(s, t)^2 + w'(s, t)^2) ds \right] u'' \\ &\quad - 2\mu_2 \dot{u} \cos \eta + \mu_1 \cos \zeta \cos \omega_0 t, \\ \ddot{w} &= -R^2 w'''' - P_0 w'' + \left[ \int_0^\pi (u'(s, t)^2 + w'(s, t)^2) ds \right] w'' \\ &\quad - 2\mu_2 \dot{w} \sin \eta + \mu_1 \sin \zeta \cos \omega_0 t, \\ u(0, t) &= u(\pi, t) = u''(0, t) = u''(\pi, t) \\ &= w(0, t) = w(\pi, t) = w''(0, t) = w''(\pi, t) = 0, \end{aligned} \quad (13)$$

where  $R^2$  is constant representing the stiffness ratio for the two directions. We assume  $R > 1$  which amounts to choosing  $w$  as the direction with stiffer cross-section. Note that  $R = 1$  reduces to Section 4.2. As before we assume  $\eta, \zeta \in (0, \pi/2)$  and  $R^2 < P_0 < 4$ . So  $1 < R < 2$ . Then we define

$$\begin{aligned} a_1^2 &= P_0 - 1, & \omega_{n-1,1}^2 &= n^2[(n^2 - P_0)], & n &= 2, 3, \dots; \\ a_2^2 &= P_0 - R^2, & \omega_{n-1,2}^2 &= n^2[n^2 R^2 - P_0], & n &= 2, 3, \dots. \end{aligned}$$

Now, the Melnikov function (9) is

$$M(\alpha) = \left[ \frac{8\omega_0 \cos \zeta}{\sqrt{\pi}} \sin \omega_0 \alpha \operatorname{sech} \frac{\pi\omega_0}{2a_1} \right] \mu_1 - \left( \frac{16a_1^3 \cos \eta}{3\pi} \right) \mu_2.$$

The following result is obtained from Theorem 3.1 (cf. [7]).

**Theorem 4.3** *If  $\omega_0 \neq \omega_{n,i}$  for all  $n$  and for  $i = 1, 2$ , then whenever  $\mu_1$  and  $\mu_2 \neq 0$  are small satisfying one of the following conditions*

$$|\mu_2| < \frac{3\sqrt{\pi} \omega_0 \cos \zeta}{2a_1^3 \cos \eta} \operatorname{sech} \frac{\pi\omega_0}{2a_1} \cdot |\mu_1|, \quad |\mu_2| < \frac{3\sqrt{\pi} \omega_0 \sin \zeta}{2a_2^3 \sin \eta} \operatorname{sech} \frac{\pi\omega_0}{2a_2} \cdot |\mu_1|,$$

*PDE (13) has a homoclinic solution with the associated chaos.*

In the  $\mu_1$ - $\mu_2$  plane in this case we get a diagram as in Figure 3. For parameter values in the regions labeled *I* there is one homoclinic orbit while for those in *II* there are two.

### 4.4 Multiple Buckled Modes

It remains to consider the situation where the axial load,  $P_0$ , is increased sufficiently to produce multiple buckled modes. So we suppose there exists an integer  $N \in \mathbb{N}$  such that  $N^2 < P_0 <$

$(N+1)^2$ . We then define

$$\begin{aligned} a_n^2 &= n^2(P_0 - n^2), \quad \text{for } n = 1, 2, \dots, N; \\ \omega_{n-N}^2 &= n^2(n^2 - P_0), \quad \text{for } n = N+1, N+2, \dots \end{aligned}$$

The following result is obtained from Theorem 3.1 (cf. [7]).

**Theorem 4.4** *Let  $N \in \mathbb{N}$ ,  $N^2 < P_0 < (N+1)^2$  and suppose one of the following hold:*

(i)  *$N$  is odd and set  $m = N$ .*

(ii)  *$N$  is even,  $N \geq 4$ , set  $m = N-1$  and*

$$P_0 \neq \frac{4N^2 - (N-1)^2 [\sqrt{9N^2 - 2N + 1} - 3(N-1)]^2}{4N^2 - [\sqrt{9N^2 - 2N + 1} - 3(N-1)]^2}.$$

(iii)  *$N = 2$ , set  $m = 1$  and*

$$P_0 \neq \frac{37 + 5\sqrt{33}}{16}, \quad P_0 \neq \frac{55 + 9\sqrt{33}}{16}.$$

*Suppose in addition that  $\omega_n \neq \omega_0$  for all  $n$ . Then whenever  $\mu_1$  and  $\mu_2 \neq 0$  are small satisfying*

$$|\mu_2| < \frac{3m\sqrt{\pi}\omega_0}{2a_m^3} \operatorname{sech} \frac{\pi\omega_0}{2a_m} \cdot |\mu_1|,$$

*PDE (1) has a homoclinic solution with the associated chaos.*

We look at the case of a beam constrained to planer motion. The calculations for the non-planer case are similar.

## Acknowledgement

The paper was supported by the Grant VEGA-MS 1/2001/05 and by the Slovak Research and Development Agency under the contract No. APVV-0414-07.

## References

- [1] BATTELLI, F., FEČKAN, M.: *Chaos in the beam equation*, In J. Differential Equations, Vol. 209, pp. 172-227, 2005.
- [2] BATTELLI, F., FEČKAN, M., FRANCA, M.: *On the chaotic behavior of a compressed beam*, In Dynamics PDE, Vol. 4, pp. 55-86, 2007.
- [3] BATTELLI, F., LAZARRI, C.: *Exponential dichotomies, heteroclinic orbits, and Melnikov functions*, In J. Differential Equations, Vol. 86., pp. 342-366, 1990.

- [4] BERTI, M., CARMINATI, C.: *Chaotic dynamics for perturbations of infinite dimensional Hamiltonian systems*, In Nonlinear Analysis, Vol. 48, pp. 481-504, 2000.
- [5] FEČKAN, M.: *Higher dimensional Melnikov mappings*, In Math. Slovaca, Vol. 49, pp. 75-83, 1999.
- [6] FEČKAN, M., GRUENDLER, J.: *The existence of chaos for ordinary differential equations with a center manifold*, In Bull. Belgian Math. Soc., Vol. 11, pp. 77-94, 2004.
- [7] FEČKAN, M., GRUENDLER, J.: *The existence of chaos in infinite dimensional non-resonant systems*, In Dynamics PDE, Vol. 5, pp. 185-209, 2008.
- [8] GRUENDLER, J.: *Homoclinic solutions for autonomous dynamical systems in arbitrary dimension*, In SIAM J. Math. Anal. Vol. 23, pp. 702-721, 1992.
- [9] HOLMES, P.: *A nonlinear oscillator with a strange attractor*, In Phil. Trans. Roy. Soc. A, Vol. 292, pp. 419-448, 1979.
- [10] HOLMES, P., Marsden, J.: *A partial differential equation with infinitely many periodic orbits: chaotic oscillations of a forced beam*, In Arch. Rational Mech. Anal., Vol. 76, pp. 135-165, 1981.
- [11] LI, Y.: *Persistent homoclinic orbits for nonlinear Schrödinger equation under singular perturbation*, In Dynamics PDE, Vol. 1, pp. 87-123, 2004.
- [12] LI, Y.: *Smale horseshoes and symbolic dynamics in perturbed nonlinear Schrödinger equations*, In J. Nonlinear Sciences, Vol. 9, pp. 363-415, 1999.
- [13] PALMER, K. J.: *Exponential dichotomies and transversal homoclinic points*, In J. Differential Equations, Vol. 55, pp. 225-256, 1984.
- [14] ROBINSON, C.: *Dynamical Systems. Stability, Symbolic Dynamics, and Chaos*, CRC Press, Boca Raton, 1995.
- [15] YAGASAKI, K.: *Homoclinic and heteroclinic behavior in an infinite-degree-of-freedom Hamiltonian system: chaotic free vibrations of an undamped, buckled beam*, In Phys. Lett. A, Vol. 285, pp. 55-62, 1991.

## Current address

### Michal Fečkan

Department of Mathematical Analysis and Numerical Mathematics, Comenius University,  
Mlynská dolina, 842 48 Bratislava, Slovakia,  
e-mail: Michal.Feckan@fmph.uniba.sk



## SINGULAR INITIAL PROBLEM FOR FREDHOLM-VOLTERRA INTEGRODIFFERENTIAL EQUATIONS

FILIPPOVA Olga, (CZ), ŠMARDA Zdeněk, (CZ)

**Abstract.** In the paper existence and uniqueness of solutions of singular Fredholm-Volterra integrodifferential equations are studied and, moreover, conditions of continuous dependence of solutions on a parameter are determined. Solutions of given integrodifferential equations are located in cone-shaped area, which gives a bound for solutions of the investigated singular problem.

**Key words and phrases.** Fredholm-Volterra integrodifferential equations, Banach fixed point theorem.

*Mathematics Subject Classification.* 45J05.

### 1 Introduction

In the past few years, many papers are devoted to the study of singular problems for differential and integrodifferential equations (see[1-9]). The fundamental methods of investigation of all above mentioned works are based on applications of fixed point theorems especially Schauder's theorem in [6], Schauder-Tychonoff's theorem in [1], Banach fixed point theorem in [5,7].

In this paper we extend results of the paper [7] to Fredholm-Volterra integrodifferential equations

$$y'(t) = \mathcal{F} \left( t, y(t), \int_{0^+}^t K_1(t, s, y(t), y(s)) ds, \int_{0^+}^1 K_2(t, s, y(t), y(s)) ds, \mu \right),$$
$$y(0^+, \mu) = 0, \tag{1}$$

and, moreover, we shall also investigate a problem of continuous dependence of solutions on a parameter.

Suppose

- (I)  $\mathcal{F} : \Omega \rightarrow R^n$ ,  $\mathcal{F} \in C^0(\Omega)$ ,  
 $\Omega = \{(t, u_1, u_2, u_3, \mu) \in J \times (R^n)^3 \times R : |u_1| \leq \phi(t), |u_2| \leq \psi(t), |u_3| \leq \psi(t)\}$ ,  $J = (0, 1]$ ,  
 $0 < \phi(t) \in C^0(J)$ ,  $\phi(0^+) = 0$ ,  $0 < \psi(t) \in C^0(J)$ ,  $|\cdot|$  denotes the usual norm in  $R^n$ ,  
 $|\mathcal{F}(t, \bar{u}_1, \bar{u}_2, \bar{u}_3, \mu) - \mathcal{F}(t, \bar{\bar{u}}_1, \bar{\bar{u}}_2, \bar{\bar{u}}_3, \mu)| \leq \sum_{i=1}^3 M_i |\bar{u}_i - \bar{\bar{u}}_i|$  for all  
 $(t, \bar{u}_1, \bar{u}_2, \bar{u}_3, \mu), (t, \bar{\bar{u}}_1, \bar{\bar{u}}_2, \bar{\bar{u}}_3, \mu) \in \Omega$ ,  $M_i \geq 0$ ,  $i = 1, 2, 3$ .

- (II)  $K_j : \Omega^1 \rightarrow R^n$ ,  $K_j \in C^0(\Omega_1)$ ,  $\Omega_1 = \{(t, s, w, v) \in J \times J \times R^n \times R^n : |w| \leq \phi(t), |v| \leq \phi(t)\}$ ,  
 $|K_1(t, s, \bar{w}, \bar{v}) - K_1(t, s, \bar{\bar{w}}, \bar{\bar{v}})| \leq [N_1 |\bar{w} - \bar{\bar{w}}| + N_2 |\bar{v} - \bar{\bar{v}}|]$ ,  
 $|K_2(t, s, \bar{w}, \bar{v}) - K_2(t, s, \bar{\bar{w}}, \bar{\bar{v}})| \leq [N_3 e^{\lambda(s-t)} |\bar{w} - \bar{\bar{w}}| + N_4 e^{\lambda(t-s)} |\bar{v} - \bar{\bar{v}}|]$   
for all  $(t, s, \bar{w}, \bar{v}), (t, s, \bar{\bar{w}}, \bar{\bar{v}}) \in \Omega_1$ ,  $N_j \geq 0$ ,  $j = 1, 2$ .  $\lambda > 0$  is a sufficiently large constant  
such that

$$\left( \frac{M_1 + M_2 N_1 + M_3 N_3 + M_3 N_4}{\lambda} + \frac{M_2 N_2}{\lambda^2} \right) < 1.$$

## 2 Main results

**Theorem 2.1** *Let the functions  $\mathcal{F}(t, u_1, u_2, u_3, \mu)$ ,  $K_j(t, s, w, v)$ ,  $j = 1, 2$  satisfy conditions (I), (II) and, moreover*

$$|\mathcal{F}| \leq \sum_{i=1}^3 g_i(t) |u_i|, \quad 0 < g_i(t) \in C^0(J), \quad \int_{0^+}^t g_1(s) \phi(s) ds \leq \alpha \phi(t),$$

$$\int_{0^+}^t (g_2(s) + g_3(s)) \psi(s) ds \leq \beta \phi(t), \quad \alpha + \beta \leq 1,$$

then the problem (1) has a unique solution  $y(t, \mu)$  for each  $\mu \in R$ ,  $t \in J$ .

**Proof.** Denote  $H$  the Banach space of continuous vector-valued functions

$$h : J_0 \rightarrow R^n, \quad J_0 = [0, 1], \quad |h(t)| \leq \phi(t)$$

on  $J$  with the norm

$$\|h\|_\lambda = \max_{t \in J_0} \{e^{-\lambda t} |h(t)|\},$$

where  $\lambda > 0$  is an arbitrary parameter. The initial value problem (1) is equivalent to the system of integral equations

$$y(t) = \int_{0^+}^t \mathcal{F} \left( s, y(s), \int_{0^+}^s K_1(s, w, y(s), y(w)) dw, \int_{0^+}^1 K_2(s, w, y(s), y(w)) dw, \mu \right) ds \quad (2)$$

Define the operator  $T$  by right-hand side of (2)

$$T(h) = \int_{0+}^t \mathcal{F} \left( s, h(s), \int_{0+}^s K_1(s, w, h(s), h(w)) dw, \int_{0+}^1 K_2(s, w, h(s), h(w)) dw, \mu \right) ds,$$

where  $h \in H$ . Let  $\mu \in R$  be fixed. The transformation  $T$  maps  $H$  continuously into itself because

$$\begin{aligned} |T(h)| &\leq \int_{0+}^t \left| \mathcal{F} \left( s, h(s), \int_{0+}^s K_1(s, w, h(s), h(w)) dw, \int_{0+}^1 K_2(s, w, h(s), h(w)) dw, \mu \right) \right| ds \leq \\ &\leq \int_{0+}^t \left[ g_1(s)|h(s)| + g_2(s) \left| \int_{0+}^s K_1(s, w, h(s), h(w)) dw \right| + g_3(s) \left| \int_{0+}^1 K_2(s, w, h(s), h(w)) dw \right| \right] ds \leq \\ &\leq \int_{0+}^t (g_1(s)\phi(s) + g_2(s)\psi(s) + g_3(s)\psi(s)) ds \leq (\alpha + \beta)\phi(t) \leq \phi(t) \end{aligned}$$

for every  $h \in H$ .

Using (I), (II) and the definition  $\|\cdot\|_\lambda$  we have

$$\begin{aligned} |T(h_2) - T(h_1)| &\leq \\ &\leq \int_{0+}^t \left| \mathcal{F} \left( s, h_2(s), \int_{0+}^s K_1(s, w, h_2(s), h_2(w)) dw, \int_{0+}^1 K_2(s, w, h_2(s), h_2(w)) dw, \mu \right) - \right. \\ &\quad \left. - \mathcal{F} \left( s, h_1(s), \int_{0+}^s K_1(s, w, h_1(s), h_1(w)) dw, \int_{0+}^1 K_2(s, w, h_1(s), h_1(w)) dw, \mu \right) \right| ds \leq \\ &\leq \int_{0+}^t \left( M_1|h_2(s) - h_1(s)| + M_2 \int_{0+}^s |K_1(s, w, h_2(s), h_2(w)) - K_1(s, w, h_1(s), h_1(w))| dw + \right. \\ &\quad \left. + M_3 \int_{0+}^1 |K_2(s, w, h_2(s), h_2(w)) - K_2(s, w, h_1(s), h_1(w))| dw \right) ds \leq \\ &\leq \int_{0+}^t \left( M_1|h_2(s) - h_1(s)| + M_2 \int_{0+}^s [N_1|h_2(s) - h_1(s)| + N_2|h_2(w) - h_1(w)|] dw + \right. \\ &\quad \left. + M_3 \int_{0+}^1 [N_3e^{\lambda(s-w)}|h_2(s) - h_1(s)| + N_4e^{\lambda(s-w)}|h_2(w) - h_1(w)|] dw \right) ds \leq \\ &\leq \|h_2 - h_1\|_\lambda \left( M_1 \int_{0+}^t e^{\lambda s} ds + M_2 N_1 \int_{0+}^t \int_{0+}^s e^{\lambda s} dw ds + \right. \\ &\quad \left. + M_2 N_2 \int_{0+}^t \int_{0+}^s e^{\lambda w} dw ds + M_3 N_3 \int_{0+}^t \int_{0+}^1 e^{\lambda(s-w)} e^{\lambda w} dw ds + M_3 N_4 \int_{0+}^t \int_{0+}^1 e^{\lambda(s-w)} e^{\lambda w} dw ds \right) = \\ &= \|h_2 - h_1\|_\lambda \left( M_1 \left( \frac{e^{\lambda t} - 1}{\lambda} \right) + M_2 N_1 \left( \frac{te^{\lambda t}}{\lambda} - \frac{e^{\lambda t} - 1}{\lambda^2} \right) + M_2 N_2 \left( \frac{e^{\lambda t} - 1}{\lambda^2} - \frac{t}{\lambda} \right) + \right. \end{aligned}$$

$$\begin{aligned}
& + M_3 N_3 \left( \frac{e^{\lambda t} - 1}{\lambda} \right) + M_3 N_4 \left( \frac{e^{\lambda t} - 1}{\lambda} \right) \leq \\
& \leq \|h_2 - h_1\|_{\lambda} e^{\lambda t} \left( \frac{M_1 + M_2 N_1 + M_3 N_3 + M_3 N_4}{\lambda} + \frac{M_2 N_2}{\lambda^2} \right).
\end{aligned}$$

Thus

$$\|T(h_2) - T(h_1)\|_{\lambda} = \max_{t \in J_0} \{e^{-\lambda t} |T(h_2) - T(h_1)|\} \leq q \|h_2 - h_1\|_{\lambda},$$

where

$$q := \frac{M_1 + M_2 N_1 + M_3 N_3 + M_3 N_4}{\lambda} + \frac{M_2 N_2}{\lambda^2}.$$

By Banach theorem the operator  $T$  has a unique stationary point  $h^*$  in the space  $H$ , i.e.  $h^*(t) \equiv T(h^*(t))$ ,  $t \in J_0$ . Then  $y := h^*$  is the desired solution of (1).

**Theorem 2.2** *Let the assumptions of Theorem 2.1 be satisfied and let there exist a constant  $L > 0$  and the integrable function  $\gamma : J_0 \rightarrow J_0$ , such that*

$$|\mathcal{F}(t, u_1, u_2, u_3, \mu_2) - \mathcal{F}(t, u_1, u_2, u_3, \mu_1)| \leq \gamma(t) |\mu_2 - \mu_1|,$$

where  $(t, u_1, u_2, u_3, \mu_1), (t, u_1, u_2, u_3, \mu_2) \in \Omega$  and

$$\max_{t \in J_0} \left\{ e^{-\lambda t} \int_{0+}^t \gamma(s) ds \right\} \leq L,$$

then the solution  $y(t, \mu)$  of (1) is continuous with respect to the variables  $(t, \mu) \in J \times R$ .

**Proof.** Define as above, for  $h \in H$  the transformation  $T_{\mu}(h)$  by means of the right-hand side (2) then we obtain

$$\|T_{\mu}(h) - T_{\mu}(y)\|_{\lambda} \leq \left( \frac{M_1 + M_2 N_1 + M_3 N_3 + M_3 N_4}{\lambda} + \frac{M_2 N_2}{\lambda^2} \right) \|h - y\|_{\lambda}.$$

By the hypothesis of Theorem 2.2 we get

$$\begin{aligned}
& e^{-\lambda t} |T_{\mu_2}(h) - T_{\mu_1}(h)| \leq \\
& e^{-\lambda t} \int_{0+}^t \left| \mathcal{F} \left( s, h(s), \int_{0+}^s K_1(s, w, h(s), h(w)) dw, \int_{0+}^1 K_2(s, w, h(s), h(w)) dw, \mu_2 \right) - \right. \\
& \left. - \mathcal{F} \left( s, h(s), \int_{0+}^s K_1(s, w, h(s), h(w)) dw, \int_{0+}^1 K_2(s, w, h(s), h(w)) dw, \mu_1 \right) \right| ds \leq \\
& \leq e^{-\lambda t} \int_{0+}^t \gamma(s) |\mu_2 - \mu_1| ds \leq L |\mu_2 - \mu_1|.
\end{aligned}$$



Hence

$$\|T_{\mu_2}(h) - T_{\mu_1}(h)\|_{\lambda} \leq L|\mu_2 - \mu_1|.$$

From this and by Theorem 2.1 we obtain

$$\begin{aligned} \|h(t, \mu_2) - h(t, \mu_1)\|_{\lambda} &= \|T_{\mu_2}[h(t, \mu_2)] - T_{\mu_2}[h(t, \mu_1)] + T_{\mu_2}[h(t, \mu_1)] - T_{\mu_1}[h(t, \mu_1)]\|_{\lambda} \leq \\ &\|T_{\mu_2}[h(t, \mu_2)] - T_{\mu_2}[h(t, \mu_1)]\|_{\lambda} + \|T_{\mu_2}[h(t, \mu_1)] - T_{\mu_1}[h(t, \mu_1)]\|_{\lambda} \leq \\ &\leq \left( \frac{M_1 + M_2N_1 + M_3N_3 + M_3N_4}{\lambda} + \frac{M_2N_2}{\lambda^2} \right) \|h(t, \mu_2) - h(t, \mu_1)\|_{\lambda} + L|\mu_2 - \mu_1|. \end{aligned}$$

Thus

$$\|h(t, \mu_2) - h(t, \mu_1)\|_{\lambda} \leq \left[ 1 - \left( \frac{M_1 + M_2N_1 + M_3N_3 + M_3N_4}{\lambda} + \frac{M_2N_2}{\lambda^2} \right) \right]^{-1} L|\mu_2 - \mu_1|.$$

Consequently the function  $h(t, \mu)$  is uniformly continuous with respect to the variable  $\mu \in R$ ; so  $y(t, \mu)$  is also continuous with respect to two variables  $(t, \mu) \in J \times R$ . The proof is complete.

**Example.** Consider the following initial problem

$$\begin{aligned} y'(t) &= \frac{t}{3}y(t) + 2t^2 \int_{0^+}^t \sqrt{s}e^{-\frac{1}{ts}} (y(t) + 2y(s)) ds + \int_{0^+}^1 e^{10(s-t)} \arctan \frac{\mu^2}{s} \left( \frac{y(t)}{2} + y(s) \right) ds + \sqrt{t^3 + 1} \\ y(0^+, \mu) &= 0. \end{aligned} \quad (3)$$

Now we can put

$$M_1 = 1/3, \quad M_2 = 2, \quad M_3 = 1, \quad N_1 = 1/e, \quad N_2 = 2/e, \quad N_3 = \pi/4, \quad N_4 = \pi/2, \quad \lambda = 10,$$

then

$$q = \frac{M_1 + M_2N_1 + M_3N_3 + M_3N_4}{\lambda} + \frac{M_2N_2}{\lambda^2} = \frac{1/3 + 2/e + \pi/4 + \pi/2}{10} + \frac{4/e}{100} < 1.$$

Now

$$|F| \leq t/3|u_1| + 2t^2|u_2| + |u_3| \Rightarrow g_1(t) = t/3, \quad g_2(t) = 2t^2, \quad g_3(t) = 1.$$

Putting  $\phi(t) = t^5/2$ ,  $\psi(t) = t^5$  we obtain

$$\int_0^t g_1(s)\phi(s)ds = \int_0^t \frac{s^6}{6}ds = \frac{t^7}{42} \leq \frac{1}{21}\phi(t) \Rightarrow \alpha = \frac{1}{21}.$$

$$\int_0^t (g_2(s) + g_3(s))\psi(s)ds = \int_0^1 (2s^2 + 1)s^5ds \leq \frac{5}{12}t^5 \leq \frac{10}{12}\phi(t) \Rightarrow \beta = 10/12.$$

Finally, we have

$$(\alpha + \beta) = 1/21 + 10/12 < 1.$$

Now from Theorem 2.1 there exists a unique solution of initial problem(3) such that  $|y(t, \mu)| \leq \frac{t^2}{2}$ .

### **Acknowledgement**

This research has been supported by the Czech Ministry of Education in the frames of MSM002160503 Research Intention MIKROSYN New Trends in Microelectronic Systems and Nanotechnologies and MSM0021630529 Research Intention Intelligent Systems in Automation.

### **References**

- [1] BENCHOHRA, M., *On initial value problem for an integrodifferential equation*, Ann. Stiint. Univ. AL.I.CUZA IASI ,XV, (1999), 297-304.
- [2] DIBLÍK, J., RUŽIČKOVÁ, M., *Existence of positive solutions of a singular initial problem for a nonlinear system of differential equations*, Rocky Mountain Journal of Mathematics **34**, (2004), 923-944.
- [3] DIBLÍK, J., NOWAK, Ch., *A nonuniqueness criterion for a singular system of two ordinary differential equations*, Nonlinear Analysis, **64** (2006), 637-656.
- [4] O'REGAN, D., *Nonresonant nonlinear singular problems in the limit circle case*, J.Math. Anal. Appl. **197**, (1996), 708-725.
- [5] ŠMARDÁ, Z., *On solutions of an implicit singular system of integrodifferential equations depending on a parameter*, Demonstratio Mathematica, Vol.XXXI, No 1, (1998), 125-130.
- [6] ŠMARDÁ, Z., *On an initial value problem for singular integro- differential equations*, Demonstratio Mathematica, Vol. XXXV, No 4, (2002), 803-811.
- [7] ŠMARDÁ, Z., *Existence and uniqueness of solutions of nonlinear integrodifferential equations*, Journal of Applied Mathematics, Statistics and Informatics, **2**, (2005), 73-79.
- [8] YANG, G., *Minimal positive solutions to some singular second-order differential equations*, J. Math. Anal. Appl. **266** (2002), 479-491.
- [9] ZHANG, B.,G., KONG, L., *Positive solutions for a class of singular boundary value problems*, Dyn. Sys. Appl.**9** (2001), 281-289.

### **Current address**

#### **Ing. Olga Filippova**

Department of Mathematics

Faculty of Electrical Engineering and Communication

Brno University of Technology, Technická 8, 616 00 BRNO

e-mail: xfilip03@stud.feec.vutbr.cz

#### **Doc. RNDr. Zdeněk Šmarda, CSc.**

Department of Mathematics

Faculty of Electrical Engineering and Communication  
Brno University of Technology, Technická 8, 616 00 BRNO  
e-mail: smarda@feec.vutbr.cz



## OSCILLATIONS FOR A HYPERBOLIC DIFFUSION EQUATION WITH TIME-DEPENDENT COEFFICIENTS

HERRMANN Leopold, (CZ), MLS Jiří, (CZ), ONDOVČIN Tomáš, (CZ)

**Abstract.** Oscillatory properties of a hyperbolic reaction-diffusion-convection equation with time-dependent coefficients are studied.

**Key words and phrases.** Oscillations, oscillatory time, reaction-diffusion-convection equation with relaxation, modified Fourier law, modified Darcy law.

*Mathematics Subject Classification.* Primary 35B05, 35L20, 76R50; Secondary 80A20, 86A05.

### 1 Hyperbolic diffusion equation

Let  $\Omega \subset \mathbb{R}^n$  be a bounded domain with sufficiently regular boundary  $\partial\Omega$ ,  $\tau_0$  a positive constant,  $t \mapsto \alpha(t)$ ,  $\alpha: \mathbb{R}^+ \rightarrow \mathbb{R}$  a continuously differentiable function,  $t \mapsto \beta(t)$ ,  $\beta: \mathbb{R}^+ \rightarrow \mathbb{R}$  a continuous function,  $\inf_{t \in \mathbb{R}^+} \beta(t) > 0$  and  $\alpha$ ,  $\dot{\alpha}$  and  $\beta$  bounded functions.

Further, let  $A(x) = (a_{jk}(x))_{j,k=1}^n$  be a matrix of functions from  $C^1(\bar{\Omega})$ , which is symmetric ( $a_{jk}(x) = a_{kj}(x)$ ,  $x \in \bar{\Omega}$ ) and positive definite uniformly with respect to  $x \in \bar{\Omega}$ , i. e. there exists  $a_0 > 0$  such that

$$\sum_{j,k=1}^n a_{jk}(x) \xi_j \xi_k \geq a_0 \sum_{j=1}^n \xi_j^2, \quad (\xi_j)_{j=1}^n \in \mathbb{R}^n, \quad x \in \bar{\Omega}. \quad (1)$$

Let  $B(x) = (b_j(x))_{j=1}^n$  be a vector of functions from  $C^1(\bar{\Omega})$  and  $c \in C(\bar{\Omega})$ . (Operators  $\operatorname{div}$  and  $\operatorname{grad}$  act solely on  $x$  variables.)

We deal with the following evolution differential equation

$$\tau_0 \frac{\partial^2 u}{\partial t^2} + \alpha(t) \frac{\partial u}{\partial t} + \beta(t) L u = 0, \quad (t, x) \in \mathbb{R}^+ \times \Omega, \quad (2)$$

where

$$L u = -\operatorname{div}(A(x) \operatorname{grad} u) + B(x) \operatorname{grad} u + c(x) u. \quad (3)$$

Equations of this type arise, for example, in the theory of the heat conduction, diffusion theory or in the theory of the fluid flow in porous medium. If the fundamental balance law (mass, energy conservation law)

$$\frac{\partial \eta}{\partial t} + \operatorname{div} \mathbf{w} = 0 \quad (4)$$

is combined with the classical constitutive laws (Fourier, Fick, and Darcy, respectively) of the form

$$\mathbf{w} = -K \operatorname{grad} \eta + \beta \mathbf{v}, \quad (5)$$

a classical parabolic reaction-diffusion-convection equation arises. If the constitutive law is modified by adding a time derivative flux term we get

$$\tau \frac{\partial \mathbf{w}}{\partial t} + \mathbf{w} = -K \operatorname{grad} \eta + \beta \mathbf{v}, \quad (6)$$

and after introducing into (4) we obtain a hyperbolic equation of type (2). (For the fluid flow in porous medium see [19].) Here  $\eta$  is the density of the substance,  $\mathbf{w}$  is the flux-density vector,  $\mathbf{v}$  is the flow field for the medium in which the substance is moving,  $K$  is a coefficient, in general a tensor, and  $\tau$  is the so-called relaxation time.

This modification of the classical (time-independent) constitutive relation (Fourier law) in the theory of heat conduction is due to Cattaneo [4]. Since then many papers have appeared concerning the Cattaneo-type heat models, “non-Fickian” diffusion, Darcy law with relaxation, *etc.*: see e. g. [1], [3], [8], [9], [17], [19]. In the study of the general motion of the fluid flow through movable matrix instead of adopting the Cattaneo approach Mls in [16] made use of the D’Alembert principle for both phases. He obtained a system of quasi-linear first-order hyperbolic equations. These equations govern the general Darcian mechanics of two-phase systems.

## 2 Oscillatory properties

In this contribution we are interested in oscillatory solutions of Eq. (2) as such solutions are frequently important in practical problems, for example when studying the tidal effects in groundwater, see e. g. [3], [18]. We shall study properties of a function  $(t, x) \mapsto u(t, x)$ ,

$u: \mathbb{R}^+ \times \Omega \rightarrow \mathbb{R}$  that solves Eq. (2) supplemented with the homogeneous Dirichlet condition  $u = 0$  for  $(t, x) \in \mathbb{R}^+ \times \partial\Omega$ . By a solution we mean any weak solution satisfying

$$\frac{d^2}{dt^2}(u(t), w) + \alpha(t) \frac{d}{dt}(u(t), w) + \beta(t) (u(t), L^+ w) = 0 \quad (7)$$

for any  $w \in W_2^2(\Omega) \overset{\circ}{W}_2^1(\Omega)$  in the sense of distributions on  $\mathbb{R}^+$ . Here  $(\cdot, \cdot)$  is the scalar product in  $L^2(\Omega)$  and

$$L^+ w = -\operatorname{div}(A(x) \operatorname{grad} w) - \operatorname{div}(B(x) w) + c(x) w \quad (8)$$

We will assume that Eq. (2) (together with the homogenous Dirichlet condition) has global (defined for all  $t \in \mathbb{R}^+$ ) solutions of finite energy  $u \in C(\mathbb{R}^+, \overset{\circ}{W}_2^1(\Omega)) \cap C^1(\mathbb{R}^+, L_2(\Omega))$  and the solutions possess the pseudo-analyticity property: if  $u$  is a solution on  $\mathbb{R}^+ \times \Omega$ ,  $T \geq 0$ ,  $\epsilon > 0$ ,

$$u = 0 \text{ on } (T, T + \epsilon) \times \Omega \implies u \equiv 0 \text{ on } \mathbb{R}^+ \times \Omega. \quad (9)$$

Let us recall that (in accordance with [6], [12], [13], [14]) a measurable function  $u: \mathbb{R}^+ \times \Omega \rightarrow \mathbb{R}$  is said to be globally oscillatory (about zero at  $+\infty$ ) if there exists (the so-called oscillatory time)  $\Theta > 0$  such that for any interval  $J \subset \mathbb{R}^+$ , the length  $|J|$  of which is greater than  $\Theta$ , the function  $u$  changes the sign on  $J \times \Omega$ , i. e. we have simultaneously  $\operatorname{meas} \{ (t, x) \in J \times \Omega \mid u(t, x) > 0 \} > 0$  and  $\operatorname{meas} \{ (t, x) \in J \times \Omega \mid u(t, x) < 0 \} > 0$ .

For  $u$  satisfying the property (9) an equivalent definition is possible:  $u: \mathbb{R}^+ \times \Omega \rightarrow \mathbb{R}$ ,  $u \not\equiv 0$  in  $\mathbb{R}^+ \times \Omega$  is globally oscillatory if and only if there exists  $\Theta > 0$  such that for any interval  $J \subset \mathbb{R}^+$  the following implication holds

$$u \geq 0 \text{ (or } u \leq 0) \text{ on } J \times \Omega \implies |J| \leq \Theta. \quad (10)$$

Roughly speaking, this means, for a continuous function  $u (\not\equiv 0)$ , that  $u$  has a zero in any domain  $J \times \Omega$  where  $J \subset \mathbb{R}^+$  is an interval the length of which is sufficiently large and this length can be chosen independently of  $J$ .

For the operator  $L^+$  supplemented with the homogeneous Dirichlet boundary condition it is well-known (cf. [2], [7], [21]) that under some regularity assumptions on the boundary  $\partial\Omega$  and coefficients there exist the so-called principal eigenvalue  $\lambda_1$  and an associated (principal) eigenfunction  $v_1$ . This means that  $L^+ v_1 = \lambda_1 v_1$  in  $\Omega$ , both  $\lambda_1$  and  $v_1$  are *real*,  $\lambda_1$  is a simple eigenvalue, i. e.  $v_1$  spans the null space  $\ker(-L^+ + \lambda_1)$ ,  $v_1$  is positive in  $\Omega$ , if  $\psi$  is a positive eigenfunction with eigenvalue  $\lambda$ , then  $\lambda = \lambda_1$ , for any eigenvalue  $\lambda$ :  $\Re \lambda \geq \lambda_1$ . Moreover, the function  $v_1$  is bounded or, in fact,  $v_1 \in C(\bar{\Omega})$ . We shall assume  $\lambda_1 > 0$  (which is equivalent to the validity of “the maximum principle”; this happens, for example, if the function  $c$  is nonnegative). In [2] also various bounds on  $\lambda_1$ , especially positive lower bounds, are established. In the sequel we use, in particular:

$$L^+ v_1 = \lambda_1 v_1 \text{ in } \Omega, \quad \lambda_1 > 0, \quad v_1 > 0 \text{ in } \Omega, \quad v_1 = 0 \text{ on } \partial\Omega. \quad (11)$$

### 3 Main result

**Theorem:** *Let*

$$\inf_{t \in \mathbb{R}^+} \left( \frac{\lambda_1 \beta(t)}{\tau_0} - \frac{\dot{\alpha}(t)}{2\tau_0} - \frac{\alpha^2(t)}{4\tau_0^2} \right) = \omega^2 > 0. \quad (12)$$

*Then Eq. (2) (under homogeneous Dirichlet condition) is uniformly globally oscillatory, i. e. there exists  $\Theta > 0$  such that any solution is globally oscillatory with the oscillatory time  $\Theta$ ,  $\Theta$  is given by formula*

$$\Theta = \frac{\pi}{\omega}, \quad \omega = \sqrt{\inf_{t \in \mathbb{R}^+} \left( \frac{\lambda_1 \beta}{\tau_0} - \frac{\dot{\alpha}}{2\tau_0} - \frac{\alpha^2}{4\tau_0^2} \right)}. \quad (13)$$

*Proof.* Let us make the projection of the equation on  $\ker(-L^+ + \lambda_1)$  and define

$$u_1(t) = \int_{\Omega} u(t, x) v_1(x) dx. \quad (14)$$

We obtain the ordinary differential equation (where  $\cdot = d/dt$ )

$$\ddot{u}_1 + \frac{\alpha}{\tau_0} \dot{u}_1 + \frac{\beta \lambda_1}{\tau_0} u_1 = 0, \quad t \in \mathbb{R}^+. \quad (15)$$

Let us assume  $u \geq 0$  (or  $u \leq 0$ ) on  $J \times \Omega$ . Owing to (11), the positivity of  $v_1$ , we get  $u_1 \geq 0$  (or  $u_1 \leq 0$ ) on  $J$ . The results of [20] (Section 8) together with the assumption (12) give:  $|J| > \frac{\pi}{\omega} \implies u_1 \equiv 0$ . Using again the positivity of  $v_1$ , we obtain  $u \equiv 0$  on  $J \times \Omega$ . Finally, due to the property (9) we get  $u \equiv 0$  on  $\mathbb{R}^+ \times \Omega$  and this completes the proof.

*Remark.* The method of the proof is based on [15], for conservative systems used in [5]. Similar results on the hyperbolic diffusion equation with constant coefficients can be obtained by means of the method described in [13] and based on results of [11], for more details see [19].

### Acknowledgement

The research has been supported by the Research Plan MSM 6840770010 and grant of GACR No. 205/07/1311.

### References

- [1] AURIAULT, J.-L., LEWANDOWSKA, J. and ROYER, P.: *About Non-Fickian Hyperbolic Diffusion*. Proc. 18<sup>ème</sup> Congrès Français de Mécanique, 27–31 août, 2007, Grenoble.
- [2] BERESTYCKI, H., NIRENBERG, L., and VARADHAN, S. R. S.: The principal eigenvalue and maximum principle for second-order elliptic operators in general domains. *Comm. Pure Appl. Math.* **47** (1994), 47–92.



- [3] BODVARSSON, G.: Confined fluids as strain meters. *J. Geophysical Research* **75** (1970), no. 14, 2711–2718.
- [4] CATTANEO, C.: Sur une forme de l'équation de la chaleur éliminant le paradoxe d'une propagation instantanée. *Comptes Rendus Acad. Sci. Paris* **247** (1958), 431–433.
- [5] CAZENAVE, T., and HARAUX, A.: Propriétés oscillatoires des solutions de certaines équations des ondes semi-linéaires. *C.R. Acad. Sc. Paris* **298** Sér. I no. 18 (1984), 449–452.
- [6] CAZENAVE, T., and HARAUX, A.: Some oscillatory properties of the wave equation in several space dimensions. *J. Functional Analysis* **76** (1988), 87–109.
- [7] GILBARG, D., and TRUDINGER N. S.: *Elliptic partial differential equations of the second order*. 2<sup>nd</sup> ed. Grundlehren der mathematischen Wissenschaften, No. 224, Springer-Verlag, Berlin-New York 1983.
- [8] GÓMEZ, H., COLOMINAS, I., NAVARRINA, F., and CASTELEIRO, M.: A finite element formulation for a convection-diffusion equation based on Cattaneo's law. *Computer Methods in Applied Mechanics and Engineering* **196** (2007), no. 9–12, 1757–1766.
- [9] GÓMEZ, H., COLOMINAS, I., NAVARRINA, F., and CASTELEIRO, M.: A mathematical model and a numerical model for hyperbolic mass transport in compressible flows. *Heat and Mass Transfer* **45** (2008), no. 2, 219–226.
- [10] HERRMANN, L.: *Diffusivity, hyperbolicity, oscillatoricity*. Proc. Sem. Topical Problems of Fluid Mechanics 2000, Inst. Thermomechanics AS CR, 17–20.
- [11] HERRMANN, L.: Conjugate points of second order ordinary differential equations with jumping nonlinearities. *Mathematics Comput. Simulations* **76** (2007), no. 1-3, 82–85.
- [12] HERRMANN, L.: Oscillations for Liénard type equations. *J. Mathématiques Pures Appl.* **90** (2008), no. 1, 60–65.
- [13] HERRMANN, L.: Differential inequalities and equations in Banach spaces with a cone. *Nonlinear Analysis, Theory, Methods and Applications* **69** (2008), no. 1, 245–255.
- [14] HERRMANN, L.: Oscillations for evolution equations with square root operators. *J. Applied Mathematics* **1** (2008), no. 1, 159–168.
- [15] HERRMANN, L., and FIALKA, M.: Oscillatory properties of equations of mathematical physics with time-dependent coefficients. *Publicationes Mathematicæ Debrecen* **57** (2000), 79–84.
- [16] MLS, J.: A continuum approach to two-phase porous media. *Transport in Porous Media* **35** (1999), 15–36.
- [17] NAVARRINA, F., GÓMEZ, H., COLOMINAS, I., and CASTELEIRO, M.: *Recent achievements on the use of pure hyperbolic Cattaneo-type convection-diffusion models in CFD*. Proc. 8th World Congress on Computational Mechanics, June 30–July 5, 2008, Venice.
- [18] ONDOVČIN, T., HERRMANN, L., and MLS, J.: Tidal effects in groundwater (Czech.) In: *Mathematics at Universities VII, Determinism and Chaos*, ed. L. Herrmann, Union of Czech Mathematicians and Physicists, 2007, 100–105.

- [19] ONDOVČIN, T., MLS, J., and HERRMANN, L.: Oscillations for an equation arising in groundwater flow with the relaxation time. (Preprint.)
- [20] PROTTER, M. H., and WEINBERGER, H. F.: *Maximum principles in differential equations*. Prentice–Hall, Inc., Englewood Cliffs, New Jersey 1967.
- [21] RABINOWITZ, P. H.: *Théorie du degré topologique et applications à des problèmes aux limites non linéaires* (rédigé par H. Berestycki). Cours de D. E. A. Publ. Lab. d'Analyse Numérique, Université Paris VI, 1975.

### **Current addresses**

#### **Doc. RNDr. Leopold Herrmann, CSc.**

Institute of Technical Mathematics, Faculty of Mechanical Engineering,  
Czech Technical University in Prague,  
Karlovo náměstí 13, 121 35 Praha 2, Czech Republic  
e-mail address: Leopold.Herrmann@fs.cvut.cz

#### **Doc. RNDr. Jiří Mls, CSc.**

Institute of Hydrogeology, Engineering Geology and Applied Geophysics, Faculty of Science,  
Charles University, Prague,  
Albertov 6, 128 43 Praha 2, Czech Republic  
e-mail address: mls@natur.cuni.cz

#### **Mgr. Tomáš Ondovčín**

Institute of Hydrogeology, Engineering Geology and Applied Geophysics, Faculty of Science,  
Charles University, Prague,  
Albertov 6, 128 43 Praha 2, Czech Republic  
e-mail address: toon@seznam.cz

## EXISTENCE OF NONOSCILLATORY SOLUTIONS OF NONLINEAR DELAY DIFFERENTIAL EQUATIONS

ILAVSKÁ Iveta, (SK), NAJMANOVÁ Anna, (SK),  
OLACH Rudolf, (SK)

**Abstract.** The article deals with the nonlinear delay differential equations. The sufficient conditions for the existence of the nonoscillatory bounded solutions are established.

**Key words and phrases.** Nonlinear differential equation; delay; existence; nonoscillatory solution; Banach space.

*Mathematics Subject Classification.* Primary 34K15.

### 1 Introduction

In the present paper we consider the nonlinear delay differential equation of the form

$$\dot{x}(t) = p(t)x(t) - q(t)x(t)x(\tau(t)), \quad t \geq t_0, \quad (1)$$

where  $p, q \in C([t_0, \infty), [0, \infty))$ ,  $p(t) \not\equiv 0$ ,  $q(t) \not\equiv 0$ ,  $\tau \in C([t_0, \infty), (0, \infty))$  is increasing function,  $\tau(t) \leq t$  and  $\lim_{t \rightarrow \infty} \tau(t) = \infty$ .

By a solution  $x(t)$  of Eq.(1) we mean a function  $x \in C([T - \tau(T), \infty), R)$  for some  $T \geq t_0$  and such that the Eq.(1) is satisfied for  $t \geq T$ .

A solution  $x(t)$  of Eq.(1) is said to be oscillatory if it has arbitrarily large zeros; otherwise it is called nonoscillatory.

The autonomous ordinary differential equation

$$\frac{dN(t)}{dt} = rN(t) \left( 1 - \frac{N(t)}{K} \right), \quad t \geq 0, \quad (2)$$

where  $r, K \in (0, \infty)$  is known as the logistic equation in mathematical ecology. The modification of Eq.(2) has the form

$$\frac{dN(t)}{dt} = rN(t) \left( 1 - \frac{N(t - \tau)}{K} \right), \quad t \geq 0, \quad (3)$$

where  $r, \tau, K \in (0, \infty)$ . The Eq.(3) is commonly known as the delay logistic equation and represents the dynamics of a single species population model. Here  $N(t)$  denotes the density of the population at time  $t$ . The symbol  $r$  is the growth rate and  $K$  is the carrying capacity of the environment. The term  $1 - N(t - \tau)/K$  denotes a feedback mechanism which takes  $\tau$  units of time to respond to changes in the size of the population.

Our interest is focused on Eq.(1) which is a generalization of the delay logistic Eq.(3). We shall develop some sufficient conditions for the existence of nonoscillatory bounded solutions of Eq.(1).

The following fixed point lemma will be used to prove the main result in the next section.

**Lemma 1.1** [2] (*Krasnoselskii's Fixed Point Theorem*)

Let  $X$  be a Banach space, let  $\Omega$  be a bounded closed convex subset of  $X$  and let  $S_1, S_2$  be maps of  $\Omega$  into  $X$  such that  $S_1x + S_2y \in \Omega$  for every pair  $x, y \in \Omega$ . If  $S_1$  is a contractive and  $S_2$  is completely continuous then the equation

$$S_1x + S_2x = x$$

has a solution in  $\Omega$ .

## 2 Existence of nonoscillatory solutions

**Theorem 2.1** Suppose that

$$\int_{t_0}^{\infty} p(t) dt < \infty \quad (4)$$

and

$$\int_{t_0}^{\infty} q(t) dt < \infty. \quad (5)$$

Then Eq.(1) has a nonoscillatory bounded solution.

**Proof.** With regard to (4) and (5) we choose a  $T > t_0$  sufficiently large such that

$$\int_T^{\infty} p(s) ds \leq \frac{1}{3}$$

and

$$\int_T^\infty Mq(s) ds \leq c,$$

where  $c > 0$ ,  $M = \max_{c \leq x, y \leq 3c} \{xy\}$ .

By  $C([t_0, \infty), R)$  we denote the set of all continuous bounded functions with the norm

$$\|x\| = \sup_{t \geq t_0} |x(t)|.$$

Then  $C([t_0, \infty), R)$  is a Banach space.

We define a closed, bounded and convex subset  $\Omega$  of  $C([t_0, \infty), R)$  as follows

$$\Omega = \{x = x(t) \in C([t_0, \infty), R) : c \leq x(t) \leq 3c, t \geq t_0\}.$$

We now define two maps  $S_1$  and  $S_2$ :  $\Omega \rightarrow C([t_0, \infty), R)$  as follows

$$(S_1x)(t) = \begin{cases} 2c - \int_t^\infty p(s)x(s) ds, & t \geq T, \\ (S_1x)(T), & t_0 \leq t \leq T, \end{cases}$$

$$(S_2x)(t) = \begin{cases} \int_t^\infty q(s)x(s)x(\tau(s)) ds, & t \geq T, \\ (S_2x)(T), & t_0 \leq t \leq T. \end{cases}$$

We shall show that for any  $x, y \in \Omega$  we have  $S_1x + S_2y \in \Omega$ . For every  $x, y \in \Omega$  and  $t \geq T$  we get

$$\begin{aligned} (S_1x)(t) + (S_2y)(t) &= 2c - \int_t^\infty p(s)x(s) ds + \int_t^\infty q(s)y(s)y(\tau(s)) ds \\ &\leq 2c + \int_t^\infty Mq(s) ds \leq 3c. \end{aligned}$$

For  $t \in [t_0, T]$  we have

$$(S_1x)(t) + (S_2y)(t) = (S_1x)(T) + (S_2y)(T) \leq 3c.$$

Furthermore for  $t \geq T$  we obtain

$$(S_1x)(t) + (S_2y)(t) \geq 2c - \int_t^\infty p(s)x(s) ds \geq 2c - 3c \int_t^\infty p(s) ds \geq 2c - c = c.$$

For  $t \in [t_0, T]$  we have

$$(S_1x)(t) + (S_2y)(t) = (S_1x)(T) + (S_2y)(T) \geq c.$$

Thus we have proved that  $S_1x + S_2y \in \Omega$  for any  $x, y \in \Omega$ . We shall show that  $S_1$  is a contraction mapping on  $\Omega$ . For  $x, y \in \Omega$  and  $t \geq T$  we have

$$\begin{aligned} |(S_1x)(t) - (S_1y)(t)| &= \left| \int_t^\infty p(s)[x(s) - y(s)] ds \right| \leq \int_t^\infty p(s)|x(s) - y(s)| ds \\ &\leq \int_t^\infty p(s) ds \|x - y\| \leq \frac{1}{3} \|x - y\|. \end{aligned}$$

This implies that

$$\|S_1x - S_1y\| \leq \frac{1}{3} \|x - y\|.$$

Also for  $t \in [t_0, T]$  we obtain the inequality above. We conclude that  $S_1$  is a contraction mapping on  $\Omega$ .

We now show that  $S_2$  is completely continuous. First we shall show that  $S_2$  is continuous. Let  $x_k = x_k(t) \in \Omega$  be such that  $x_k(t) \rightarrow x(t)$  as  $k \rightarrow \infty$ . Because  $\Omega$  is closed,  $x = x(t) \in \Omega$ . For  $t \geq T$  we get

$$|(S_2x_k)(t) - (S_2x)(t)| \leq \int_t^\infty q(s)|x_k(s)x_k(\tau(s)) - x(s)x(\tau(s))| ds.$$

Since

$$|x_k(s)x_k(\tau(s)) - x(s)x(\tau(s))| \rightarrow 0 \quad \text{as } k \rightarrow \infty,$$

by applying the Lebesgue dominated convergence theorem we obtain that

$$\lim_{k \rightarrow \infty} \|(S_2x_k)(t) - (S_2x)(t)\| = 0.$$

This means that  $S_2$  is continuous. We now show that  $S_2\Omega$  is relatively compact. By the Arzela-Ascoli theorem it is sufficient to show that the family of functions  $\{S_2x : x \in \Omega\}$  is uniformly bounded and equicontinuous on  $[t_0, \infty)$ . The uniform boundedness follows from the definition of  $\Omega$ . For the equicontinuity we only need to show that for any given  $\epsilon > 0$  the interval  $[t_0, \infty)$  can be decomposed into finite subintervals in such a way that on each subinterval all functions of the family have change of amplitude less than  $\epsilon$ . With regard to the conditions of theorem for any  $\epsilon > 0$  we take  $T^* \geq T$  large enough so that

$$\int_{T^*}^\infty q(s)x(s)x(\tau(s)) ds \leq \frac{\epsilon}{2}.$$

Then for  $x \in \Omega$ ,  $T_2 > T_1 \geq T^*$  we get

$$\begin{aligned} |(S_2x)(T_2) - (S_2x)(T_1)| &\leq |(S_2x)(T_2)| + |(S_2x)(T_1)| \\ &= \int_{T_2}^\infty q(s)x(s)x(\tau(s)) ds + \int_{T_1}^\infty q(s)x(s)x(\tau(s)) ds \leq \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon. \end{aligned}$$

For  $x \in \Omega$  and  $T \leq T_1 < T_2 \leq T^*$  we obtain

$$|(S_2x)(T_2) - (S_2x)(T_1)| = \int_{T_1}^{T_2} q(s)x(s)x(\tau(s)) ds \leq \max_{T \leq s \leq T^*} \{q(s)x(s)x(\tau(s))\}(T_2 - T_1).$$

Then there exists  $\delta > 0$  such that

$$|(S_2x)(T_2) - (S_2x)(T_1)| < \epsilon \quad \text{if } 0 < T_2 - T_1 < \delta.$$

For any  $x \in \Omega$ ,  $t_0 \leq T_1 < T_2 \leq T$  we have

$$|(S_2x)(T_2) - (S_2x)(T_1)| = 0 < \epsilon.$$

Then  $\{S_2x : x \in \Omega\}$  is uniformly bounded and equicontinuous on  $[t_0, \infty)$ . Hence  $S_2\Omega$  is relatively compact. By Lemma 1.1 there is  $x_0 \in \Omega$  such that  $S_1x_0 + S_2x_0 = x_0$ . Thus  $x_0(t)$  is a positive bounded solution of Eq.(1). The proof is complete.

**Corollary 2.2** Assume that

$$\int_{t_0}^{\infty} p(t) dt < \infty.$$

Then nonlinear delay differential equation

$$\dot{x}(t) = p(t)[x(t) - x(t)x(\tau(t))], \quad t \geq t_0,$$

has a nonoscillatory bounded solution.

**Example 1.** Consider nonlinear delay differential equation

$$\dot{x}(t) = t \exp(-t)x(t) - \exp(-2t)x(t)x(\tau(t)), \quad t \geq 0.$$

Since conditions (4), (5) are satisfied, this equation has a nonoscillatory bounded solution.

**Example 2.** Consider nonlinear neutral delay differential equation

$$\dot{x}(t) = \exp(-t)x(t) - t \exp(-t)x(t)x(\tau(t)), \quad t \geq 0.$$

Since conditions (4), (5) are satisfied, this equation has a nonoscillatory bounded solution.

## Acknowledgement

The research was supported by the grant APVV-0700-07 and 1/0090/09 of Scientific Grant Agency of Ministry of Education of Slovak Republic.

## References

- [1] DIBLÍK, J. - KÚDELČÍKOVÁ, M.: Two classes of asymptotically different positive solutions of the equation  $\dot{y} = -f(t, y_t)$ , Nonlinear Analysis, article in press.
- [2] ERBE, L.H. - KONG, Q.K. - ZHANG, B.G.: Oscillation Theory for Functional Differential Equations. Marcel Dekker, New York, 1995.
- [3] GYÖRI, I. - LADAS, G.: Oscillation Theory of Delay Differential Equations. Clarendon Press, Oxford, 1991.
- [4] MARUŠIAK, P. - OLACH, R.: Functional Differential Equations. EDIS-ŽU, Žilina, 2000 (in Slovak language).
- [5] YU, Y. H. - WANG, H. Z.: Nonoscillatory solutions of second-order nonlinear neutral delay equations. J. Math. Anal. Appl. 331(2005), 445-456.
- [6] ZHOU, Y.: Existence for nonoscillatory solutions of second-order nonlinear differential equations. J. Math. Anal. Appl. 331(2007), 91-96.
- [7] ZHOU, Y. - ZHANG, B. G.: Existence of nonoscillatory solutions of neutral differential equations with positive and negative coefficients. Appl. Math. Lett. 15(2002), 867-874.

## Current address

### **ILAVSKÁ, IVETA, Mgr.**

University of Žilina, Faculty of Science  
Department of Mathematical Analysis and Applied Mathematics  
J.M.Hurbana St.15,  
010 26 Žilina, Slovakia  
**e-mail:** *iveta.ilavska@fpv.uniza.sk*

### **NAJMANOVÁ, ANNA, Mgr.**

University of Žilina, Faculty of Science  
Department of Mathematical Analysis and Applied Mathematics  
J.M.Hurbana St.15,  
010 26 Žilina, Slovakia  
**e-mail:** *anna.najmanova@fpv.uniza.sk*

### **OLACH, RUDOLF, Doc., RNDr., CSc.**

University of Žilina, Faculty of Science  
Department of Mathematical Analysis and Applied Mathematics  
J.M.Hurbana St.15 ,  
010 26 Žilina, Slovakia  
**e-mail:** *rudolf.olach@fpv.uniza.sk*



## THE $\theta$ -METHODS FOR THE DELAY DIFFERENTIAL EQUATIONS

JÁNSKÝ Jiří, (CZ)

**Abstract.** This paper deals with the discretization of the delay differential equations. We pay the special attention to the  $\theta$ -methods for these equations. In particular, we study the qualitative behaviour of solutions of this method applied to the pantograph equation.

**Key words and phrases.** The  $\theta$ -method, pantograph equation, asymptotic estimate.

*Mathematics Subject Classification.* Primary 39A11, 39A12.

### 1 Introduction

We consider the differential equation with a delayed argument in the form

$$y'(t) = a(t)y(t) + b(t)y(\phi(t)), \quad t \geq t_0, \quad (1)$$

where  $a(t), b(t), \phi(t)$  are continuous and real function on  $[t_0, \infty)$  and  $\phi(t) < t$  for  $t > t_0$ ,  $\phi(t_0) \leq t_0$ . The popular discretization of the equation (1) is the well-known  $\theta$  method involving e.g. Euler methods and trapezoidal rule as particular cases. The different approaches to this type of a discretization are mentioned in [1, 11]. The aim of this paper is twofold: First we describe and distinguish this approaches and comment relations between them. Secondly, we consider the pantograph equation as the particular case of (1) via the choose  $\phi(t) = \lambda t$ ,  $0 < \lambda < 1$ ,  $a(t) = a$ ,  $b(t) = b$ ,  $t_0 = 0$  and describe the asymptotic behaviour of the  $\theta$ -methods applied to the pantograph equation. We note, that this and related problems have been extensively studied, e.g. in [2 - 10, 12]. On this account we make also some comparisons with the known results.

## 2 The derivation of the $\theta$ -methods

Integration of (1) yields

$$\int_0^t y'(\tau) d\tau = \int_0^t a(\tau) y(\tau) d\tau + \int_0^t b(\tau) y(\phi(\tau)) d\tau. \quad (2)$$

We introduce the substitution  $u = \phi(\tau)$  and denote

$$\psi(u) := \phi^{-1}(u).$$

Then equation (2) becomes

$$y(t) - y(0) = \int_0^t a(\tau) y(\tau) d\tau + \int_0^{\phi(t)} b(\phi^{-1}(\tau)) \psi'(\tau) y(\tau) d\tau.$$

After the discretization we get

$$y_{n+1} - y_n = \int_{t_0+nh}^{t_0+(n+1)h} a(\tau) y(\tau) d\tau + \int_{\phi(t_0+nh)}^{\phi(t_0+(n+1)h)} b(\phi^{-1}(\tau)) \psi'(\tau) y(\tau) d\tau. \quad (3)$$

The integrals on the right-hand side of (3) can be approximated by use of the explicit rectangular formula as well as implicit rectangular formula. We denote:  $y_n \approx y(t_0+nh)$ ,  $b_n = b(t_0+nh)$ ,  $a_n = a(t_0+nh)$  and  $\phi_n = \phi(t_0+nh)$ .

First we approximate both integrals on the right-hand side of the equation (3) using the rectangular formula with the left grid point, i.e.

$$\int_{t_0+nh}^{t_0+(n+1)h} a(\tau) y(\tau) d\tau \approx h a_n y_n,$$

$$\int_{\phi_n}^{\phi_{n+1}} b(\phi^{-1}(\tau)) \psi'(\tau) y(\tau) d\tau \approx (\phi_{n+1} - \phi_n) b_n \psi'(\phi_n) y(\phi_n).$$

The equation (3) becomes

$$y_{n+1} = y_n + h a_n y_n + b_n (\phi_{n+1} - \phi_n) \psi'(\phi_n) y(\phi_n). \quad (4)$$

Since the point  $\phi_n$  is not usually a grid point, we define the value  $y(\phi_n)$  as the linear interpolation

$$y(\phi_n) := (1 - r_n) y_{\lfloor \frac{\phi_n - t_0}{h} \rfloor} + r_n y_{\lfloor \frac{\phi_n - t_0}{h} \rfloor + 1}, \quad (5)$$

where  $r_n := \frac{\phi_n - t_0}{h} - \lfloor \frac{\phi_n - t_0}{h} \rfloor$ .

Now we proceed to another way of discretization, which is based on the fact that integrals on the right-hand side of the equation (3) are approximated using the rectangular formulae with the right grid point. Since the substitution of the first integral is quite simple, it is omitted here. The substitution of the second integral has the form

$$\int_{\phi_n}^{\phi_{n+1}} b(\phi^{-1}(\tau))\psi'(\tau)y(\tau)d\tau \approx (\phi_{n+1} - \phi_n)b_{n+1}\psi'(\phi_{n+1})y(\phi_{n+1}).$$

After performing all the operations mentioned above we get

$$y_{n+1} = y_n + ha_{n+1}y_{n+1} + b_{n+1}(\phi_{n+1} - \phi_n)\psi'(\phi_{n+1})y(\phi_{n+1}). \quad (6)$$

Now we have the same problem as above. The point  $\phi_{n+1}$  is not usually a grid point. Thus we define the value  $y(\phi_{n+1})$  as the linear interpolation. To simplify the resulting relation we use the same grid points as in (5) i.e.

$$y(\phi_{n+1}) := (1 - k_n)y_{\lfloor \frac{\phi_n - t_0}{h} \rfloor} + k_n y_{\lfloor \frac{\phi_n - t_0}{h} \rfloor + 1}, \quad (7)$$

where  $k_n := \frac{\phi_{n+1} - t_0}{h} - \lfloor \frac{\phi_n - t_0}{h} \rfloor$ . We note that the value  $k_n$  can be greater than 1.

The linear combination of (4) and (6) yields

$$\begin{aligned} y_{n+1} &= y_n + h((1 - \theta)a_n y_n + \theta a_{n+1} y_{n+1}) \\ &\quad + (\phi_{n+1} - \phi_n)((1 - \theta)b_n \psi'(\phi_n) y(\phi_n) + \theta b_{n+1} \psi'(\phi_{n+1}) y(\phi_{n+1})), \end{aligned} \quad (8)$$

where  $\theta \in [0, 1]$  and  $y(\phi_n), y(\phi_{n+1})$  are given by (5) and (7). Note, that equation (8) for  $\theta = 1/2$  was derived using the procedure stated in [1].

Now we present another way of discretization of (1). Rewrite the equation (1) as

$$y_{n+1} - y_n = \int_{t_0 + nh}^{t_0 + (n+1)h} a(\tau)y(\tau)d\tau + \int_{t_0 + nh}^{t_0 + (n+1)h} b(\tau)y(\phi(\tau))d\tau. \quad (9)$$

Both integral on the right-hand side of the equation (9) are replaced as follows: They are approximated by using rectangular formulae by using left grid point at first. The approximation of the second integral is in the form

$$\int_{t_0 + nh}^{t_0 + (n+1)h} b(\tau)y(\phi(\tau))d\tau \approx hb_n y(\phi_n).$$

Thus we get

$$y_{n+1} = y_n + ha_n y_n + hb_n y(\phi_n). \quad (10)$$

Similarly we can arrive at

$$y_{n+1} = y_n + ha_{n+1} y_{n+1} + hb_{n+1} y(\phi_{n+1}). \quad (11)$$

The linear combination of (10) and (11) yields the  $\theta$  method in the form

$$y_{n+1} = y_n + h((1 - \theta)a_n y_n + \theta a_{n+1} y_{n+1}) + h((1 - \theta)b_n y(\phi_n) + \theta b_{n+1} y(\phi_{n+1})), \quad (12)$$

where  $y(\phi_n)$  and  $y(\phi_{n+1})$  are given by (5), (7) respectively. Note that this equation can be also found in [11], where  $y(\phi_{n+1})$  is calculated via the linear interpolation utilizing the left and right neighborhoods of  $\phi_{n+1}$ .

### 3 Some auxiliary results

In this section we deal with the case, where  $a(t) = a$ ,  $b(t) = b$  and  $\phi(t) = \lambda t$ ,  $0 < \lambda < 1$ ,  $t_0 = 0$ . In such a case the equation (1) becomes

$$y'(t) = ay(t) + by(\lambda t), \quad t \geq 0, \quad (13)$$

and the formula (8) as well as formula (12) give the recurrence relation

$$y_{n+1} = Ry_n + S(\beta_n y_{[\lambda n]} + \alpha_n y_{[\lambda n]+1}), \quad (14)$$

where  $R := \frac{1+(1-\theta)ah}{1-\theta ah}$ ,  $S := \frac{bh}{1-\theta ah}$ ,  $\beta_n := 1 - \alpha_n$  and  $\alpha_n := \lambda_n - [\lambda n] + \theta\lambda$ .

Furthermore we assume

$$\eta = \eta(\theta, \lambda) := \sup_{n \in \mathbb{Z}^+} (|\beta_n| + |\alpha_n|) < \infty, \quad |R| < 1.$$

Next lemma can be found in the particular case  $\theta = 1/2$  in [1, Theorem 6].

**Lemma 3.1** *Let  $0 < \lambda < 1$ ,  $0 \leq \theta \leq 1$ . Then the function  $\eta(\theta, \lambda)$  has the following values:*

$$\eta(\theta, \lambda) = \begin{cases} 1, & \lambda = K/L, \quad \theta K \leq 1, \quad K, L \in \{1, 2, 3, \dots\} \text{ and relatively prime,} \\ 1 + 2\theta\lambda - \frac{2}{L}, & \lambda = K/L, \quad \theta K \geq 1, \quad K, L \in \{2, 3, \dots\} \text{ and relatively prime,} \\ 1 + 2\theta\lambda, & \lambda \text{ irrational.} \end{cases} \quad (15)$$

**Proof.** First note that  $1 \leq \eta(\theta, \lambda) \leq 1 + 2\theta\lambda$ . Now assume  $\lambda = \frac{K}{L}$  where  $1 \leq K < L$  and  $(K, L)$  are relatively prime. It is known that

$$\frac{nK}{L} - \lfloor \frac{nK}{L} \rfloor = \frac{nK \bmod L}{L}.$$

Then

$$\sup_{n \in \mathbb{Z}^+} \alpha_n = \theta\lambda + \sup_{n \in \mathbb{Z}^+} (\lambda n - [\lambda n]) = \theta\lambda + \frac{L-1}{L} = 1 + \theta\lambda - \frac{1}{L}.$$

Thus the first two cases of (15) hold.

Let  $\lambda$  be irrational number. The case  $\theta = 0$  is trivial, hence we deal only with the case  $\theta \neq 0$ . In this case for every  $\epsilon > 0$ ,  $\epsilon < \theta\lambda$  there exist an  $n_\epsilon$  such that

$$1 - \epsilon < \lambda n_\epsilon - [\lambda n_\epsilon].$$

Furthermore,

$$\alpha_{n_\epsilon} > 1 + \theta\lambda - \epsilon > 1$$

and we arrive at

$$\eta(\theta, \lambda) \geq \alpha_{n_\epsilon} + |1 - \alpha_{n_\epsilon}| = 1 + 2\theta\lambda - 2\epsilon.$$

Now we get  $\eta(\theta, \lambda) \geq 1 + 2\theta\lambda$ , because of  $\epsilon > 0$  can be made arbitrary small.  $\square$

Now we present the inequality which is useful in our further calculations. This inequality has the form:

$$|S| (|\beta_n| \rho_{\lfloor \lambda n \rfloor} + \alpha_n \rho_{\lfloor \lambda n \rfloor + 1}) \leq (1 - |R|) \rho_n, \quad n = 0, 1, \dots \quad (16)$$

**Lemma 3.2** *The sequence*

$$\rho_n := \begin{cases} \left(n - \frac{1}{1-\lambda}\right)^{-\log_\lambda \gamma} & \text{for } \gamma \geq 1, \\ \left(n + \frac{1}{1-\lambda}\right)^{-\log_\lambda \gamma} & \text{for } 0 < \gamma < 1 \end{cases} \quad (17)$$

where

$$\gamma := \frac{|S|\eta}{1 - |R|} \quad (18)$$

defines the solution of inequality (16).

**Proof.** We only deal with the case  $0 < \gamma < 1$  because the case  $\gamma \geq 1$  is analogous. If  $0 < \gamma < 1$ , then  $(\rho_n)$  is a decreasing sequence. Hence, we can write

$$|S| (|\beta_n| \rho_{\lfloor \lambda n \rfloor} + \alpha_n \rho_{\lfloor \lambda n \rfloor + 1}) \leq |S| (|\beta_n| + \alpha_n) \rho_{\lfloor \lambda n \rfloor} = |S| \eta \rho_{\lfloor \lambda n \rfloor}.$$

Further

$$\begin{aligned} |S| \eta \rho_{\lfloor \lambda n \rfloor} &= |S| \eta \left( \lfloor \lambda n \rfloor + \frac{1}{1-\lambda} \right)^{-\log_\lambda \gamma} \\ &\leq |S| \eta \left( \lambda n - 1 + \frac{1}{1-\lambda} \right)^{-\log_\lambda \gamma} \\ &= (1 - |R|) \rho_n. \end{aligned}$$

## 4 Main result

This section presents the main result of this paper. First we introduce the following notation. If  $y_n$  is a solution of (14) and  $\rho_n$  is given by (17), then we denote

$$B_0 := \sup(|y_n|/\rho_n, n \in [\lfloor \lambda \sigma_0 \rfloor, \sigma_0] \cap \mathbb{Z}^+), \quad (19)$$

where

$$\sigma_0 \geq \max \left( \frac{2}{(1-\lambda)\lambda}, 2 \log_\lambda \gamma \right) \quad (20)$$

is arbitrary integer number. Furthermore, we denote

$$L := \max \left( \frac{\log_{\lambda} \gamma}{\gamma(1 - |R|)(\sigma_0 - \frac{1+\lambda}{1-\lambda})}, \frac{2 \log_{\lambda} \gamma}{\sigma_0 - \frac{1+\lambda}{1-\lambda}} \right). \quad (21)$$

Now we can formulate the main theorem of this paper.

**Theorem 4.1** *Let  $y_n$  be a solution of (14), where  $a < 0$ ,  $b \neq 0$ ,  $\lambda \in (0, 1)$  and let  $\gamma$ ,  $B_0$ ,  $\sigma_0$ ,  $L$  be given by (18)-(21). Then*

$$|y_n| \leq B_0 e^{\frac{L}{1-\lambda}} n^{-\log_{\lambda} \gamma} \quad \text{for } n = \sigma_0, \sigma_0 + 1, \sigma_0 + 2, \dots \quad (22)$$

**Proof.** We use the substitution  $z_n = y_n/\rho_n$  in (14), where  $\rho_n$  is given by (17). Then

$$\varrho_{n+1} z_{n+1} = R \varrho_n z_n + S \left( \beta_n \rho_{\lfloor \lambda n \rfloor} z_{\lfloor \lambda n \rfloor} + \alpha_n \rho_{\lfloor \lambda n \rfloor + 1} z_{\lfloor \lambda n \rfloor + 1} \right). \quad (23)$$

Now we choose  $\sigma_0 \geq \max(\frac{2}{(1-\lambda)\lambda}, 2 \log_{\lambda} \gamma)$ ,  $\sigma_0 \in \mathbb{Z}^+$  and define points  $\sigma_{m+1} := \lfloor \frac{\sigma_m - 1}{\lambda} \rfloor$ , where  $m = 0, 1, \dots$ . After some calculations, we obtain

$$\lambda^{-m} \left( \sigma_0 - \frac{1+\lambda}{1-\lambda} \right) \leq \sigma_m \leq \lambda^{-1} \sigma_{m-1}, \quad m = 1, 2, \dots \quad (24)$$

Next we introduce intervals  $I_0 := [\lfloor \lambda \sigma_0 \rfloor, \sigma_0] \cap \mathbb{Z}^+$ ,  $I_{m+1} := [\sigma_m, \sigma_{m+1}] \cap \mathbb{Z}^+$  and denote  $B_m := \sup(|z_s|, s \in \cup_{j=0}^m I_j)$ ,  $m = 0, 1, 2, \dots$ .

Now we choose  $n^* \in I_{m+1}$ ,  $n^* > \sigma_m$  arbitrarily and we distinguish two cases with respect to  $R$ .

(i) First, we deal with the case  $R = 0$ . In this case

$$z_{n^*} = \frac{S}{\rho_{n^*}} \left( \beta_{n^*-1} \rho_{\lfloor \lambda(n^*-1) \rfloor} z_{\lfloor \lambda(n^*-1) \rfloor} + \alpha_{n^*-1} \rho_{\lfloor \lambda(n^*-1) \rfloor + 1} z_{\lfloor \lambda(n^*-1) \rfloor + 1} \right),$$

Thus

$$|z_{n^*}| \leq B_m \frac{|S|}{\rho_{n^*}} \left( |\beta_{n^*-1}| \rho_{\lfloor \lambda(n^*-1) \rfloor} + \alpha_{n^*-1} \rho_{\lfloor \lambda(n^*-1) \rfloor + 1} \right)$$

Using (16), we get

$$|z_{n^*}| \leq (1 - |R|) \frac{\rho_{n^*-1}}{\rho_{n^*}} B_m \leq \frac{\rho_{n^*-1}}{\rho_{n^*}} B_m.$$

Assuming  $\gamma \geq 1$ ,  $(\varrho_n)$  is the nondecreasing sequence and we obtain  $|z_{n^*}| \leq B_m$ . Assuming  $0 < \gamma < 1$  we derive with respect to (17), (24) and the binomial formula the relation

$$\frac{\rho_{n^*-1}}{\rho_{n^*}} = \left( \frac{n^* + \frac{1}{1-\lambda} - 1}{n^* + \frac{1}{1-\lambda}} \right)^{-\log_{\lambda} \gamma} \leq \frac{1}{\left(1 + \frac{1}{\sigma_m}\right)^{-\log_{\lambda} \gamma}} \leq \frac{1}{1 + \frac{-\log_{\lambda} \gamma}{\sigma_m}} \leq 1 + \frac{2 \log_{\lambda} \gamma}{\sigma_m}.$$

This inequality implies the following relation

$$|z_{n^*}| \leq B_m \left( 1 + \frac{2 \log_{\lambda} \gamma}{\sigma_0 - \frac{1+\lambda}{1-\lambda}} \lambda^m \right). \quad (25)$$

(ii) Let  $R \neq 0$ . We can multiply the equation (23) by  $\frac{1}{R^{n+1}}$ . We get

$$\Delta \left( \frac{\varrho_n z_n}{R^n} \right) = \frac{S}{R^{n+1}} (\beta_n \rho_{[\lambda_n]} z_{[\lambda_n]} + \alpha_n \rho_{[\lambda_n]+1} z_{[\lambda_n]+1}).$$

If we sum this relation from  $\sigma_m$  to  $n^* - 1$ , then we obtain

$$\frac{\varrho_{n^*} z_{n^*}}{R^{n^*}} - \frac{\varrho_{\sigma_m} z_{\sigma_m}}{R^{\sigma_m}} = \sum_{p=\sigma_m}^{n^*-1} \frac{S}{R^{p+1}} (\beta_p \rho_{[\lambda_p]} z_{[\lambda_p]} + \alpha_p \rho_{[\lambda_p]+1} z_{[\lambda_p]+1}),$$

i.e.

$$z_{n^*} = \frac{\varrho_{\sigma_m}}{\varrho_{n^*}} \frac{R^{n^*}}{R^{\sigma_m}} z_{\sigma_m} + \frac{R^{n^*}}{\varrho_{n^*}} \sum_{p=\sigma_m}^{n^*-1} \frac{S}{R^{p+1}} (\beta_p \rho_{[\lambda_p]} z_{[\lambda_p]} + \alpha_p \rho_{[\lambda_p]+1} z_{[\lambda_p]+1}).$$

Thus

$$\begin{aligned} |z_{n^*}| &\leq \frac{\varrho_{\sigma_m}}{\varrho_{n^*}} \frac{|R|^{n^*}}{|R|^{\sigma_m}} |z_{\sigma_m}| + \frac{|R|^{n^*}}{\varrho_{n^*}} \sum_{p=\sigma_m}^{n^*-1} \frac{|S|}{|R|^{p+1}} |\beta_p \rho_{[\lambda_p]} z_{[\lambda_p]} + \alpha_p \rho_{[\lambda_p]+1} z_{[\lambda_p]+1}| \\ &\leq B_m \left( \frac{\varrho_{\sigma_m}}{\varrho_{n^*}} \frac{|R|^{n^*}}{|R|^{\sigma_m}} + \frac{|R|^{n^*}}{\varrho_{n^*}} \sum_{p=\sigma_m}^{n^*-1} \frac{|S|}{|R|^{p+1}} (|\beta_p| \rho_{[\lambda_p]} + \alpha_p \rho_{[\lambda_p]+1}) \right). \end{aligned}$$

Using (16), we get

$$|z_{n^*}| \leq B_m \left( \frac{\varrho_{\sigma_m}}{\varrho_{n^*}} \frac{|R|^{n^*}}{|R|^{\sigma_m}} + \frac{|R|^{n^*}}{\varrho_{n^*}} \sum_{p=\sigma_m}^{n^*-1} \frac{1 - |R|}{|R|^{p+1}} \rho_p \right).$$

Now using the relation

$$\frac{1 - |R|}{|R|^{p+1}} = \Delta \left( \frac{1}{|R|} \right)^p \quad (26)$$

and summing by parts we obtain

$$\begin{aligned} |z_{n^*}| &\leq B_m \left( \frac{\varrho_{\sigma_m}}{\varrho_{n^*}} \frac{|R|^{n^*}}{|R|^{\sigma_m}} + \frac{|R|^{n^*}}{\varrho_{n^*}} \sum_{p=\sigma_m}^{n^*-1} \Delta \left( \frac{1}{|R|} \right)^p \rho_p \right) \\ &= B_m \left( 1 - \frac{|R|^{n^*}}{\varrho_{n^*}} \sum_{p=\sigma_m}^{n^*-1} \frac{1}{|R|^{p+1}} \Delta \rho_p \right). \end{aligned}$$

Now with respect to (26), we get

$$|z_{n^*}| \leq B_m \left( 1 - \frac{|R|^{n^*}}{\varrho_{n^*}} \sum_{p=\sigma_m}^{n^*-1} \frac{\Delta \rho_p}{1 - |R|} \Delta \left( \frac{1}{|R|} \right)^p \right).$$

If  $\gamma \geq 1$  then  $\rho_p$  is nondecreasing, therefore  $\Delta\rho_p \geq 0$  and  $|z_{n^*}| \leq B_m$ . In the case  $0 < \gamma < 1$ , the same simple calculations are necessary to derive that  $\Delta\rho_p$  is negative and nondecreasing. Hence we can write

$$\begin{aligned} |z_{n^*}| &\leq B_m \left( 1 - \frac{|R|^{n^*}}{1 - |R|} \frac{\Delta\rho_{\sigma_m}}{\varrho_{n^*}} \sum_{p=\sigma_m}^{n^*-1} \Delta \left( \frac{1}{|R|} \right)^p \right) \\ &= B_m \left( 1 - \frac{|R|^{n^*}}{1 - |R|} \frac{\Delta\rho_{\sigma_m}}{\varrho_{n^*}} \left( \frac{1}{|R|^{n^*}} - \frac{1}{|R|^{\sigma_m}} \right) \right) \\ &\leq B_m \left( 1 + \frac{1}{1 - |R|} \frac{-\Delta\rho_{\sigma_m}}{\varrho_{\sigma_{m+1}}} \right). \end{aligned}$$

Substituting the corresponding form of  $\rho_n$  and using the binomial formula, we can derive

$$\begin{aligned} -\Delta\rho_{\sigma_m} &= \left( \sigma_m + \frac{1}{1 - \lambda} \right)^{-\log_\lambda \gamma} \left( 1 - \left( 1 + \frac{1}{\sigma_m + \frac{1}{1 - \lambda}} \right)^{-\log_\lambda \gamma} \right) \\ &\leq \left( \sigma_m + \frac{1}{1 - \lambda} \right)^{-\log_\lambda \gamma} \left( 1 - \left( 1 + \frac{-\log_\lambda \gamma}{\sigma_m + \frac{1}{1 - \lambda}} \right) \right) \\ &\leq \left( \sigma_m + \frac{1}{1 - \lambda} \right)^{-\log_\lambda \gamma} \frac{\log_\lambda \gamma}{\sigma_m}. \end{aligned}$$

Analogically,

$$\begin{aligned} \rho_{\sigma_{m+1}} &= \left( \sigma_{m+1} + \frac{1}{1 - \lambda} \right)^{-\log_\lambda \gamma} \\ &\geq \left( \frac{1}{\lambda} \sigma_m + \frac{1}{1 - \lambda} \right)^{-\log_\lambda \gamma} \\ &\geq \left( \frac{1}{\lambda} \sigma_m + \frac{1}{\lambda} \frac{1}{1 - \lambda} \right)^{-\log_\lambda \gamma} \\ &= \gamma \left( \sigma_m + \frac{1}{1 - \lambda} \right)^{-\log_\lambda \gamma}. \end{aligned}$$

Now we arrive at the estimate

$$\frac{-\Delta\rho_{\sigma_m}}{\rho_{\sigma_{m+1}}(1 - \tilde{R})} \leq \frac{\log_\lambda \gamma}{\gamma(1 - |R|)} \frac{1}{\sigma_m} \leq \frac{\log_\lambda \gamma}{\gamma(1 - |R|)} \frac{1}{(\sigma_0 - \frac{1+\lambda}{1-\lambda})} \lambda^m.$$

by use of (24). Hence

$$|z_{n^*}| \leq B_m \left( 1 + \frac{\log_\lambda \gamma}{\gamma(1 - |R|)(\sigma_0 - \frac{1+\lambda}{1-\lambda})} \lambda^m \right). \quad (27)$$

Summarizing cases (i)-(ii) and using estimates (25) and (27) we get

$$|z_{n^*}| \leq B_m(1 + L\lambda^m) \quad \text{as } m \rightarrow \infty$$

for arbitrary  $n^* \in I_{m+1}$ ,  $n^* > \sigma_m$ . Thus

$$B_{m+1} \leq B_m(1 + L\lambda^m) \quad \text{as } m \rightarrow \infty$$



Now we can estimate  $B_m$  in this way:

$$B_{m+1} \leq B_m (1 + L(\lambda^m)) \leq B_0 \prod_{j=0}^m (1 + L(\lambda^j)) \leq B_0 e^{\frac{L}{1-\lambda}}.$$

Thus

$$B_m \leq B_0 e^{\frac{L}{1-\lambda}} \quad \text{as } m \rightarrow \infty$$

and the estimate (22) is proved.  $\square$

**Example 4.2** Let us consider the equation (13) in the form

$$y'(t) = -y(t) - 0.5y(3t/4), \quad t \geq 0, \quad y(0) = 1. \quad (28)$$

The formula (14) with  $\theta = 1/3$  and the stepsize  $h = 0.05$  gets

$$\begin{aligned} y_0 &= 1, \\ y_{n+1} &= \frac{58}{61}y_n - \frac{30}{61}(\beta_n y_{\lfloor 3n/4 \rfloor} + \alpha_n y_{\lfloor 3n/4 \rfloor + 1}), \end{aligned}$$

where

$$\alpha_n = \frac{3n}{4} - \lfloor \frac{3n}{4} \rfloor + \frac{1}{4} < 1, \quad \beta_n = 1 - \alpha_n,$$

Now if we set  $\sigma_0 = 1000$ , then using Theorem 4.1 we obtain the estimate

$$|y_n| \leq 42n^{-2.4}, \quad \text{for } n = 1000, 1001, \dots \quad (29)$$

Note, that this estimate corresponds with estimate for exact solution (28) (see. [5, 6]).

The Fig. 1 displays the real numerical solution of the problem (28) and its estimate given by (29).

## 5 Some comparisons

In this section we compare the asymptotic estimate (22) with the known results. Let us assume first that

$$|R| + |S|\eta \leq 1, \quad a < 0, \quad \theta = 1/2.$$

Under this assumption it is shown in [1] that any solution  $y_n$  of (14) is bounded and, furthermore, the following asymptotic estimate holds:

$$y_n = O(n^{-\log_\lambda(|R|+|S|\eta)}) \quad \text{as } n \rightarrow \infty. \quad (30)$$

Let us compare now the asymptotic estimate (22) with (30). It is easy to show, that if  $|R| + |S|\eta < 1$  then

$$\gamma = \frac{|S|\eta}{1 - |R|} < |R| + |S|\eta.$$

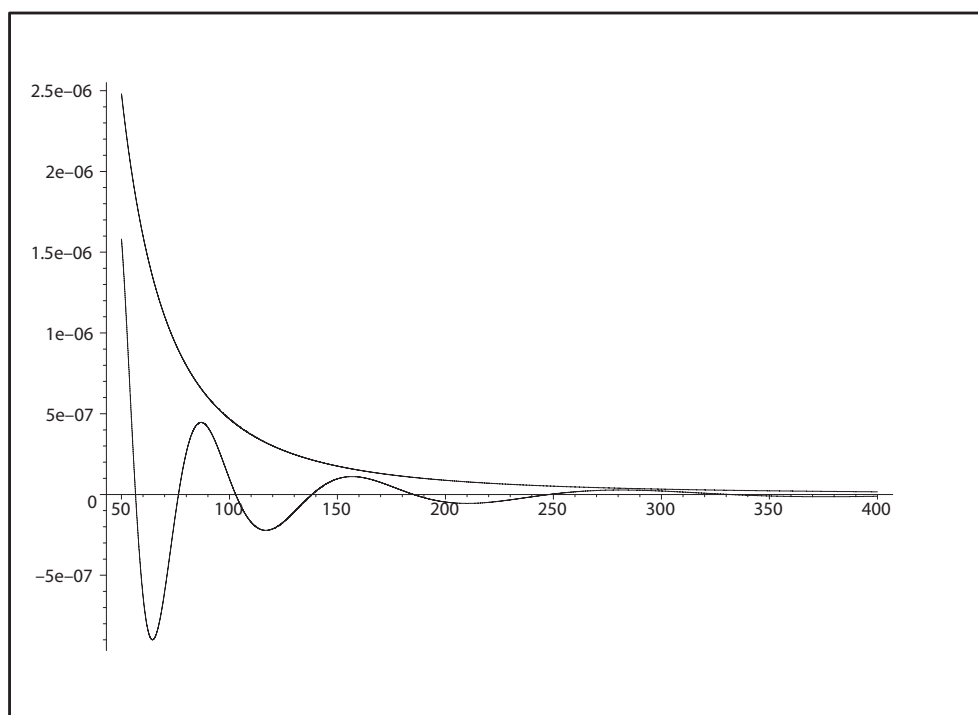


Fig. 1:

After substitution for  $|R|$  a  $|S|$  we get

$$\gamma = \begin{cases} \eta|b/a| & \text{for } (1 - \theta)h|a| \leq 1, \\ h|b|\eta/(2 + h|a|(2\theta - 1)) & \text{for } (1 - \theta)h|a| > 1. \end{cases}$$

The Lemma 3.1 implies that if  $\lambda = K/L$  where  $K, L \in \{1, 2, 3, \dots\}$  are relatively prime then we can put  $\theta = 1/K$  and obtain  $\gamma = |b/a|$  with a restriction on the stepsize  $h$ . In the case  $\theta = 1$  and  $\lambda = 1/L$  we get the value  $|b/a|$  without any restriction to the stepsize  $h$ . If  $\theta = 0$  then  $\eta = 1$  and we get the value  $|b/a|$  provided  $h \leq 1/|a|$ .

### Acknowledgement

The author was supported by the research plan MSM 0021630518 "Simulation modelling of mechatronic systems" of the Ministry of Education, Youth and Sports of the Czech Republic and by the grant # 201/08/0469 of the Czech Grant Agency.

### References

- [1] BUHMANN M.D., ISERLES A.: *Stability of the discretized pantograph differential equation*. Math. Comp., **60**, pp. 575-589, 1993.
- [2] ČERMÁK J.: *On a linear differential equation with a proportional delay*. Math. Nachr., pp. 495-504, 2007.

- [3] ČERMÁK J.: *The asymptotic of solutions for a class delay differential equations*, Rocky Mountain J. math. **33**, pp. 775-786, 2003.
- [4] ISERLES A.: *Numerical analysis of delay differential equations with variable delays*, Ann. Numer. Math. **1**, pp. 133-152, 1994.
- [5] ISERLES A.: *On the generalized pantograph functional-differential equation*, European J. Appl. Math. **4**, pp. 1-38, 1993.
- [6] KATO T., MCLEOD J.B.: *The functional-differential equation  $y'(x) = ay(\lambda x) + by(x)$* , Bull. Amer. Math. Soc. **77**, pp. 891-937, 1971.
- [7] KOTO T.: *Stability of Runge-Kutta methods for the generalized pantograph equation*, Numer. Math. **84**, pp. 870-884, 1999.
- [8] KUNDRÁT P.: *On the asymptotics of the difference equation with a proportional delay*, Opuscula Math. 26/3, pp. 499-506, 2006.
- [9] LEHNINGER H., LIU Y.: *The functional-differential equation  $y'(t) = Ay(t) + By(\lambda t) + Cy'(qt) + f(t)$* , European J. Appl. Math. **9**, pp. 81-91, 1998.
- [10] LIU Y.: *Numerical investigation of the pantograph equation*, Appl. Numer. Math. **24**, pp. 309-317, 1997.
- [11] LIU Y.: *On the  $\theta$ -method for delay differential equations with infinite lag*, J. Comput. Appl. Math. **71**, pp. 177-190, 1996.
- [12] LIU M. Z., YANG Z. W., XU Y.: *The stability of the modified Runge-Kutta methods for the pantograph equation*, Math. Comp. **75**, pp. 1201-1215, 2006.

#### Current address

**Jiří Janský, Ing.**

Brno University of Technology, Technická 2896/2, 616 69 Brno, Czech Republic,  
e-mail: yjansk04stud.fme.vutbr.cz



## A PSEUDOHYPERBOLIC PROBLEM FOR VON A KÁRMÁN SYSTEM

KEČKEMÉTYOVÁ Mária, (SK), BOCK Igor, (SK)

**Abstract.** The existence and the uniqueness of solutions is proved for dynamic problems of perpendicular vibrations of von Kármán plates whose viscosity has the character of a short memory. The boundary conditions describe the plate partly clamped and partly free. The existence of a weak solution is verified after transforming the system to one pseudohyperbolic initial-boundary value problem and using the Galerkin method. The energy dissipation is verified.

**Key words and phrases.** von Kármán plates, short memory, existence and uniqueness of solutions, energy behavior.

*Mathematics Subject Classification.* Primary 35L70, 74D10; Secondary 74K20.

### 1 Introduction and notation

Dynamic problems for viscoelastic structures represent an important but complex topic of applied mathematics. The effect of viscoelasticity for the damping of the energy of the vibrating structure plays very important role. The aim of the present paper is to study the nonlinear models of von Kármán plates. The presented results also extend the research made for the quasistatic contact problems for these plates. The case of anisotropic plate was studied in [1]. The paper [2] describes the isotropic case using the Rothe's method with respect to the time variable. In both cases the existence of solution could be proved only for sufficiently small right-hand sides. In the dynamic case the existence and the uniqueness can be verified for right-hand sides without any upper bounds. The dynamic problems for viscoelastic von Kármán plates with a long memory was studied in [8], where a viscosity memory term appears only in the equation for the deflection of the plate. We study the case, where the short memory viscosity appears in both equations of the von Kármán system. The existence of a solution is proved using the Galerkin method.

Let  $\Omega$  be a bounded domain with a Lipschitz boundary  $\Gamma = \bar{\Gamma}_0 \cup \bar{\Gamma}_1$ ,  $\Gamma_0 \cap \Gamma_1 = \emptyset$ . We assume that  $meas(\Gamma_0) > 0$  and  $\Gamma_0$  is not a straight line. The unit outer normal vector is denoted by  $\mathbf{n} = (n_1, n_2)$ ,  $\boldsymbol{\tau} = (-n_2, n_1)$  is the unit tangent vector. The displacement is denoted by  $\mathbf{u} \equiv (u_i)$ . Further we denote:

$$\frac{\partial}{\partial s} \equiv \partial_s, \quad \frac{\partial^2}{\partial s \partial r} \equiv \partial_{sr}, \quad \partial_i = \partial_{x_i}, \quad i = 1, \dots, N.$$

Let  $I \equiv (0, T)$  be a bounded time interval and

$$Q = I \times \Omega, \quad S = I \times \Gamma, \quad S_i = I \times \Gamma_i, \quad i = 0, 1.$$

We denote by  $W_p^k(M)$  with  $k \geq 0$  and  $p \in [1, \infty]$  the Sobolev spaces of functions defined on a domain or an appropriate manifold  $M$ . By  $\dot{W}_p^k(M)$  we denote the spaces with zero traces on  $\partial M$ . If  $p = 2$  we use the notation  $H^k(M)$ ,  $\dot{H}^k(M)$ . For the anisotropic spaces  $W_p^k(M)$   $k = (k_1, k_2) \in \mathbb{R}_+^2$ ,  $k_1$  is related with the time while  $k_2$  with the space variables (with the obvious consequences for  $p = 2$ ) provided  $M$  is a time-space domain. The duals to  $\dot{H}^k(M)$  are denoted by  $H^{-k}(M)$ . By  $C$  we denote the space of continuous functions with the appropriate sup-norm. For a Banach space  $X$  the space of functions  $f : I \mapsto X$  with a norm  $\|f(\cdot)\|_X \in L_p(I)$  is denoted by  $L_p(I; X)$ .

Strain tensor

$$\varepsilon_{ij}(\mathbf{u}) = \frac{1}{2}(\partial_i u_j + \partial_j u_i + \partial_i u_3 \partial_j u_3) - x_3 \partial_{ij} u_3, \quad i, j = 1, 2, \varepsilon_{i3} \equiv 0, \quad i = 1, 2, 3$$

with nonlinearities appeared in partial derivatives of perpendicular deflections corresponds to plates with moderately large deflections.

Let  $\delta_{ij}$  be the Kronecker symbol. We employ the Einstein summation convention. The viscoelastic constitutional law has the form

$$\begin{aligned} \sigma_{ij}(\mathbf{u}) = & \frac{E_1}{1 - \nu^2} \partial_t ((1 - \nu) \varepsilon_{ij}(\mathbf{u}) + \nu \delta_{ij} \varepsilon_{kk}(\mathbf{u})) \\ & + \frac{E_0}{1 - \nu^2} ((1 - \nu) \varepsilon_{ij}(\mathbf{u}) + \nu \delta_{ij} \varepsilon_{kk}(\mathbf{u})). \end{aligned}$$

The constants  $E_0, E_1 > 0$  and  $\nu \in (0, \frac{1}{2})$  are the Young modulus of elasticity, the modulus of viscosity and the Poisson ratio, respectively. In the sequel we shall use the following abbreviations:

$$a = \frac{h^2}{12}, \quad b = \frac{h^2}{12\rho(1 - \nu^2)},$$

where  $h$  is the plate thickness and  $\rho$  is the density of the material. Moreover, we introduce

$$[u, v] \equiv \partial_{11} u \partial_{22} v + \partial_{22} u \partial_{11} v - 2 \partial_{12} u \partial_{12} v, \quad u, v \in H^2(\Omega).$$

the important Poisson bracket.

The following generalization of the Aubin's compactness lemma verified in [5] Theorem 3.1 will be essentially used:

**Lemma 1.1** *Let  $B_0 \hookrightarrow\hookrightarrow B \hookrightarrow B_1$  be Banach spaces, the first reflexive and separable. Let  $1 < p < \infty$ ,  $1 \leq q < \infty$ . Then*

$$W \equiv \{v; v \in L_p(I; B_0), \dot{v} \in L_q(I, B_1)\} \hookrightarrow\hookrightarrow L_p(I; B).$$

## 2 Existence and uniqueness of a solution

### 2.1 Problem formulation

Applying the approach used in the classical theory of von Kármán equations for elastic plates (see [3]) we obtain in the dynamic case the classical formulation composed of the system

$$\left. \begin{aligned} \ddot{u} - a\Delta\ddot{u} + b(E_1\Delta^2\dot{u} + E_0\Delta^2u) - [u, v] &= f, \\ \Delta^2v + E_1\partial_t[u, u] + E_0[u, u] &= 0 \end{aligned} \right\} \text{ on } Q, \quad (1)$$

the boundary conditions

$$\begin{aligned} u = \partial_n u &= 0 \text{ on } S_0, \quad \mathcal{M}(u) = \Sigma(u) = 0 \text{ on } S_1, \\ v = \partial_n v &= 0 \text{ on } S \end{aligned} \quad (2)$$

where

$$\begin{aligned} \mathcal{M}(u) &= b[E_1M(\dot{u}) + E_0M(u)], \\ M(u) &= \Delta u + (1 - \nu)(2n_1n_2\partial_{12}u - n_1^2\partial_{22}u - n_2^2\partial_{11}u); \\ \Sigma(u) &= b[E_1V(\dot{u}) + E_0V(u)] - a\partial_n\ddot{u}, \\ V(u) &= \partial_n\Delta u + (1 - \nu)\partial_\tau[(n_1^2 - n_2^2)\partial_{12}u + n_1n_2(\partial_{22}u - \partial_{11}u)], \end{aligned}$$

and the initial conditions

$$u(0, \cdot) = u_0, \quad \dot{u}(0, \cdot) = u_1 \text{ on } \Omega. \quad (3)$$

The unknown functions  $u, v$  express the deflection of the middle plane of the plate and the Airy stress function respectively. The plate is acting upon a perpendicular load  $f$ .

For  $u, y \in L_2(I; H^2(\Omega))$  we define the following bilinear form

$$A : (u, y) \mapsto b(\partial_{kk}u\partial_{kk}y + \nu(\partial_{11}u\partial_{22}y + \partial_{22}u\partial_{11}y) + 2(1 - \nu)\partial_{12}u\partial_{12}y) \quad (4)$$

almost everywhere on  $Q$  and introduce a set

$$V = \{y \in H^2(\Omega) : y = \partial_n y = 0 \text{ on } \Gamma_0\}.$$

The set  $V$  is a Hilbert space with a scalar product and a norm

$$((y, z)) = \int_{\Omega} \Delta y \Delta z \, dx, \quad \|y\| = ((y, y))^{1/2}$$

equivalent with the obvious norm  $\|\cdot\|_2$  in a Sobolev space  $H^2(\Omega)$  (see [7], chapter 10).

The problem (1), (2), (3) has then the variational formulation:

Look for  $\{u, v\} \in H^{1,2}(Q) \times L_2(I; \dot{H}^2(\Omega))$  such that  $\dot{u} \in L_2(I; V)$ ,  $\ddot{u} \in L_2(I; H^1(\Omega))$ , the following system

$$\int_{\Omega} (a\nabla\ddot{u} \cdot \nabla z_1 + \ddot{u}z_1 + E_1A(\dot{u}, z_1) + E_0A(u, z_1) - [u, v]z_1) \, dx = \int_{\Omega} f z_1 \, dx, \quad (5)$$

$$\int_{\Omega} (\Delta v \Delta z_2 + (E_1 \partial_t [u, u] + E_0 [u, u]) z_2) dx = 0 \quad (6)$$

is satisfied for any  $(z_1, z_2) \in V \times \dot{H}^2(\Omega)$  and the conditions (3) remain valid.

We define the bilinear operator  $\Phi : H^2(\Omega)^2 \rightarrow \dot{H}^2(\Omega)$  by means of the variational equation

$$\int_{\Omega} \Delta \Phi(u, v) \Delta \varphi dx = \int_{\Omega} [u, v] \varphi dx, \quad \varphi \in \dot{H}^2(\Omega). \quad (7)$$

The equation (7) has a unique solution, because  $[u, v] \in L_1(\Omega) \hookrightarrow H^2(\Omega)^*$ . The well-defined operator  $\Phi$  is evidently compact and symmetric. Moreover, due to Lemma 1 from [6]  $\Phi : H^2(\Omega)^2 \rightarrow W_p^2(\Omega)$ ,  $2 < p < \infty$  and

$$\|\Phi(u, v)\|_{W_p^2(\Omega)} \leq c \|u\|_{H^2(\Omega)} \|v\|_{W_p^1(\Omega)} \quad \forall u \in H^2(\Omega), v \in W_p^1(\Omega). \quad (8)$$

With the help of the operator  $\Phi$  we get the following reformulation of (5,6):

**Problem  $\mathcal{P}$ .**

We look for  $u \in H^{1,2}(Q)$  such that  $\dot{u} \in L_2(I; H^2; V)$ ,  $\ddot{u} \in L_2(I; H^1(\Omega))$ , the equation

$$\begin{aligned} & \int_{\Omega} (a \nabla \ddot{u} \cdot \nabla z + \ddot{u} z + E_1 A(\dot{u}, z) + E_0 A(u, z) + [u, E_1 \partial_t \Phi(u, u) + E_0 \Phi(u, u)] z) dx \\ & = \int_{\Omega} f z dx \end{aligned} \quad (9)$$

holds for any  $z \in V$  and the conditions (3) remain valid.

We shall verify the existence and the uniqueness of a solution to the Problem  $\mathcal{P}$ .

**Theorem 2.1** *Let  $f \in L_2(Q)$ ,  $u_i \in H^2(\Omega)$ ,  $i = 0, 1$ . Then there exists a unique solution  $u \in H^{1,2}(Q)$  of the problem  $\mathcal{P}$ .*

*If  $v = -E_1 \partial_t \Phi(u_m, u_m) - E_0 \Phi(u_m, u_m)$ , then a couple  $\{u, v\}$  is a unique solution of the problem (5), (6), (3).*

*Proof. (i) The existence.* We shall apply the Galerkin method. Let us denote by  $\{w_i \in H^2(\Omega); i \in \mathbb{N}\}$  a orthonormal basis of  $H^2(\Omega)$ . We construct the Galerkin approximation  $u_m$  of a solution in a form

$$u_m(t) = \sum_{i=1}^m \alpha_i(t) w_i, \quad \alpha_i(t) \in \mathbb{R}, \quad i = 1, \dots, m, \quad m \in \mathbb{N}$$

given by the solution of the approximated problem

$$\begin{aligned} & \int_{\Omega} (a \nabla \ddot{u}_m(t) \cdot \nabla w_i + \ddot{u}_m(t) w_i + E_1 A(\dot{u}_m(t), w_i) + E_0 A(u_m(t), w_i) \\ & + [u_m(t), w_i] (E_1 \partial_t \Phi(u_m, u_m)(t) + E_0 \Phi(u_m, u_m)(t))) dx \end{aligned} \quad (10)$$

$$\begin{aligned} & = \int_{\Omega} f(t) w_i dx, \quad i = 1, \dots, m, \\ & u_m(0) = u_{0m}, \quad \dot{u}_m(0) = u_{1m}, \quad u_{im} \rightarrow u_i \text{ in } H^2(\Omega), \quad i = 0, 1. \end{aligned} \quad (11)$$



The matrix  $\mathbb{A} = (a_{ij})$ ,  $a_{ij} = \int_{\Omega} (a \nabla w_i \cdot \nabla w_j + w_i w_j) dx$  is positively definite. The system (10) can then be expressed in the form

$$\ddot{\alpha}_i = F_i(t, \dot{\alpha}_1, \dots, \dot{\alpha}_m, \alpha_1, \dots, \alpha_m), \quad i = 1, \dots, m.$$

Its right-hand side satisfies the conditions for the local existence of a solution fulfilling the initial conditions corresponding the functions  $u_{0m}, u_{1m}$ . Hence there exists a Galerkin approximation  $u_m(t)$  defined on some interval  $I_m \equiv [0, t_m]$ ,  $0 < t_m < T$ . In order to receive the conditions for the prolongation of a solution to the whole interval  $I \equiv [0, T]$  we derive the *a priori* estimates of  $\{u_m\}$  not dependent on  $t_m$ . Let  $Q_m = [0, t_m] \times \Omega$ . After multiplying the equation (10) by  $\dot{\alpha}_i(t)$ , summing up with respect to  $i$  and integrating and taking in mind the property

$$\int_{\Omega} [u, v] y dx = \int_{\Omega} [u, y] v dx, \quad (12)$$

if at least one element of  $\{u, v, y\}$  belongs to  $\dot{H}^2(\Omega)$ , cf. [3] we get

$$\begin{aligned} & \int_{Q_m} \frac{1}{2} \partial_t (a |\nabla \dot{u}_m(t)|^2 + \dot{u}_m^2(t) + E_0 A(u_m(t), u_m(t)) + \frac{E_0}{2} (\Delta \Phi(u_m, u_m))^2) + \\ & + E_1 A(\dot{u}_m(t), \dot{u}_m(t)) + \frac{E_1}{2} (\Delta \partial_t \Phi(u_m, u_m)(t))^2 dx dt = \int_{Q_m} f(t) \dot{u}_m(t) dx dt, \end{aligned} \quad (13)$$

and

$$\begin{aligned} & \|\dot{u}_m\|_{L_2(I_m; H^2(\Omega))}^2 + \|\dot{u}_m\|_{L_{\infty}(I_m; H^1(\Omega))}^2 + \|u_m\|_{L_{\infty}(I_m; H^2(\Omega))}^2 + \|\partial_t \Phi(u_m, u_m)\|_{L_2(I_m; H^2(\Omega))}^2 \\ & \leq c \equiv c(f, u_0, u_1), \end{aligned} \quad (14)$$

$$\|\partial_t \Phi(u_m, u_m)\|_{L_2(I_m; W_p^2(\Omega))} \leq c_p \equiv c_p(f, u_0, u_1) \quad \forall p > 2, \quad I_m = (0, t_m). \quad (15)$$

As the right-hand side of the estimate (14) does not depend on  $m$  we can set  $t_m = T$ .

The estimate (15) further implies

$$\begin{aligned} & [u_m, E_1 \partial_t \Phi(u_m, u_m) + E_0 \Phi(u_m, u_m)] \in L_2(I; L_r(\Omega)), \quad r = \frac{2p}{p+2}, \\ & \|[u_m, E_1 \partial_t \Phi(u_m, u_m) + E_0 \Phi(u_m, u_m)]\|_{L_2(I; L_r(\Omega))} \leq c_r \equiv c_r(f, u_0, u_1). \end{aligned} \quad (16)$$

After multiplying the equation (10) by  $\ddot{\alpha}_i(t)$ , summing up with respect to  $i$  and integrating we obtain the estimate of  $\ddot{u}$

$$\|\ddot{u}_m\|_{L_2(I; H^1(\Omega))}^2 \leq c, \quad m \in \mathbb{N}. \quad (17)$$

The estimate (16) and the imbedding  $L_q(\Omega) \hookrightarrow H^1(\Omega)$ ,  $q = \frac{2p}{p-2}$  play the crucial role in deriving (17).

We proceed with the convergence of the Galerkin approximation. Applying the estimates (14-17) and the Aubin-Lions compactness lemma we obtain a subsequence of  $\{u_m\}$  (again

denoted by  $\{u_m\}$  and a function  $u$  such that

$$\begin{aligned} \dot{u}_m &\rightharpoonup^* \dot{u} && \text{in } L_\infty(I; H^1(\Omega)), \\ \dot{u}_m &\rightharpoonup \dot{u} && \text{in } L_2(I; H^2(\Omega)), \\ \ddot{u}_m &\rightharpoonup \ddot{u} && \text{in } L_2(I; H^1(\Omega)), \\ \dot{u}_m &\rightarrow \dot{u} && \text{in } L_p(I; H^1(\Omega)) \cap L_2(I; H^{2-\varepsilon}(\Omega)) \quad \forall \varepsilon \in (0, 1), \\ u_m &\rightarrow u && \text{in } C_0(I; W_p^1(\Omega)), \\ \partial_t \Phi(u_m, u_m) &\rightharpoonup \partial_t \Phi(u, u) && \text{in } L_2(I; W_p^2(\Omega)). \end{aligned} \quad (18)$$

The fourth convergence follows for  $p \leq 2$  from the second and third one via the compact imbedding theorem, for  $p > 2$  by interpolation of the previous result with the first convergence. The fifth convergence  $u_m \rightarrow u$  in (18) follows from the second and the fourth one by a standard interpolation and imbedding technique (cf. [4], Chapter 2). The last convergence is a consequence of (8) and the second and fifth convergence.

Let  $\mu \in \mathbb{N}$  and  $z_\mu = \sum_{i=1}^m \phi_i(t) w_i$ ,  $\phi_i \in \mathcal{D}(0, T)$ ,  $i = 1, \dots, \mu$ . We have for arbitrary  $m \in \mathbb{N}$  and  $t \in I$  the relation

$$\begin{aligned} &\int_{\Omega} (a \nabla \ddot{u}_m(t) \cdot \nabla z_\mu(t) + \ddot{u}_m(t) z_\mu(t) + E_1 A(\dot{u}_m(t), z_\mu) + E_0 A(u_m(t), z_\mu(t)) \\ &+ [u_m(t), z_\mu(t)] (E_1 \partial_t \Phi(u_m, u_m)(t) + E_0 \Phi(u_m, u_m)(t))) dx = \int_{\Omega} f(t) z_\mu(t) dx. \end{aligned}$$

The convergence process (18) and the property (12) imply that a function  $u$  fulfils

$$\begin{aligned} &\int_Q (a \nabla \ddot{u} \cdot \nabla z_\mu + \ddot{u} z_\mu + E_1 A(\dot{u}, z_\mu) + E_0 A(u, z_\mu) + [u, E_1 \partial_t \Phi(u, u) + E_0 \Phi(u, u)] z_\mu) dx dt \\ &= \int_Q f z dx dt. \end{aligned}$$

Functions  $\{z_\mu\}$  form the dense subset of the set  $L_2(I; H^2(\Omega))$  and hence a function  $u$  fulfils the identity (9). The initial conditions (3) follow due to (11) and the proof of the existence of a solution is complete.

(ii) *The uniqueness.* Let  $u, \hat{u}$  be two solutions of Problem  $\mathcal{P}$  and let  $w = u - \hat{u}$ . We have for arbitrary  $s \in I$ ,  $Q_s = [0, s] \times \Omega$  the relation

$$\begin{aligned} &\int_{Q_s} (a \nabla \dot{w} \cdot \nabla z + \dot{w} z + E_1 A(\dot{w}, z) + E_0 A(w, z) \\ &+ ([u, E_1 \partial_t \Phi(u, u) + E_0 \Phi(u, u)] - [\hat{u}, E_1 \partial_t \Phi(\hat{u}, \hat{u}) + E_0 \Phi(\hat{u}, \hat{u})]) z) dx dt = 0, \\ &\forall z \in L_2(I; H^2(\Omega)), \\ &w(0, \cdot) = \dot{w}(0, \cdot) = 0 \text{ on } \Omega. \end{aligned}$$

After setting  $z = \dot{w}$  we get

$$\begin{aligned} &\frac{1}{2} (a \|\nabla \dot{w}\|_{L_2(\Omega)^2} + \|\dot{w}\|_{L_2(\Omega)} + E_0 \int_{\Omega} A(w, w) dx)(s) + E_1 \int_{Q_s} A(\dot{w}, \dot{w}) dx dt \\ &= \int_{Q_s} ([\hat{u}, E_1 \partial_t \Phi(\hat{u}, \hat{u}) + E_0 \Phi(\hat{u}, \hat{u})] - [u, E_1 \partial_t \Phi(u, u) + E_0 \Phi(u, u)]) \dot{w} dx dt. \end{aligned}$$

Using the estimate (16) with  $u$  and  $\hat{u}$  instead of  $u_m$  and the imbedding  $L_q(\Omega) \hookrightarrow H^1(\Omega)$ ,  $q = \frac{2p}{p-2}$  we obtain the inequality

$$\|\dot{w}\|_{H^1(\Omega)}^2(s) \leq c \int_0^s \|\dot{w}\|_{H^1(\Omega)}^2(t) dt \text{ for every } s \in I.$$

The Gronwall lemma implies

$$\|\dot{w}\|_{H^1(\Omega)}(s) = 0 \text{ for every } s \in I$$

and the uniqueness of a solution follows due to zero initial conditions for  $w \equiv u - \hat{u}$ .

### 3 Behavior of the energy

The aim of this section is to analyze the asymptotic behavior of the complete energy of the plate as  $t$  tends to  $\infty$ . For every  $t \geq 0$  we define the energy functional

$$\mathcal{E}(t) = \frac{1}{2} \int_{\Omega} \left( \dot{u}^2 + a|\nabla \dot{u}|^2 + E_0 A(u, u) + \frac{1}{2} E_0 (\Delta \Phi(u, u))^2 - 2fu \right) (t, x) dx. \quad (19)$$

It can be expressed as a sum of the kinetic and the potential energy of the plate acting under the force  $f$ . For simplicity we consider further the stationary right-hand e.g.  $f(t, x) \equiv f(x)$ . A following theorem expresses the decay of the energy.

**Theorem 3.1** *It holds*

$$\mathcal{E}'(t) \leq 0 \quad \forall t \in (0, \infty) \quad (20)$$

and hence

$$\mathcal{E}(t_2) \leq \mathcal{E}(t_1) \quad \forall t_1 \leq t_2, \quad t_i \in (0, \infty), \quad i = 1, 2. \quad (21)$$

Moreover, if  $f = 0$  then

$$\lim_{t \rightarrow \infty} \mathcal{E}_1(t) = 0, \quad (22)$$

where

$$\mathcal{E}_1(t) = \frac{1}{2} \int_{\Omega} (\dot{u}^2 + a|\nabla \dot{u}|^2) dx$$

corresponds to the kinetic energy of the plate.

*Proof.* After inserting  $z = \dot{u}$  in (9) we obtain the relation

$$\begin{aligned} & \frac{1}{2} \frac{d}{dt} \int_{\Omega} \left( \dot{u}^2 + a|\nabla \dot{u}|^2 + E_0 A(u, u) + \frac{1}{2} E_0 (\Delta \Phi(u, u))^2 - 2fu \right) (t, x) dx \\ &= \mathcal{E}'(t) = - \int_{\Omega} (E_1 A(\dot{u}, \dot{u}) + 2E_1 (\Delta \Phi(\dot{u}, u))^2(t, x)) dx \leq 0 \end{aligned}$$

and the decay of energy is verified.

If  $f = 0$ , then we have the relation  $\mathcal{E}(t) = \mathcal{E}(0) - \int_0^t \mathcal{E}_2(s) ds$ , where

$$\mathcal{E}_2(t) = \int_{\Omega} (E_1 A(\dot{u}, \dot{u}) + 2E_1 (\Delta \Phi(\dot{u}, u))^2(t, x)) dx.$$

We have  $\int_0^\infty \mathcal{E}_2(t) dt < \infty$  and hence  $\int_0^\infty \mathcal{E}_1(t) dt < \infty$ . Due to the uniform continuity of  $\mathcal{E}_1$  on  $(0, \infty)$  the convergence (22) follows.

**Remark 3.2** *It is possible to verify for sufficiently small right-hand sides  $f(t, x)$ ,  $t > 0$ ,  $x \in \Omega$  that in the case*

$$\lim_{t \rightarrow \infty} \int_{\Omega} (f(t, x) - f_{\infty}(x))^2 dx = 0$$

*the solution  $\{u(t, \cdot), v(t, \cdot)\}$  tends to a weak solution  $\{u_{\infty}, v_{\infty}\} \in V \times \mathring{H}^2(\Omega)$  of the corresponding elastic problem*

$$\left. \begin{aligned} E_0 \Delta^2 u - [u, v] &= f_{\infty}, \\ \Delta^2 v + E_0 [u, u] &= 0 \end{aligned} \right\} \text{ on } \Omega,$$

$$\begin{aligned} u = \partial_n u &= 0 \text{ on } \Gamma_0, \quad M(u) = V(u) = 0 \text{ on } \Gamma_1, \\ v = \partial_n v &= 0 \text{ on } \Gamma. \end{aligned}$$

**Remark 3.3** *Similar results can be obtained in the case of a long memory of a decreasing exponential type using for instance in the paper [8], where the memory appeared only in the first equation of the system.*

## Acknowledgement

The authors gratefully acknowledge the Scientific Grant Agency VEGA for supporting this work under the Grant No. 1/4214/07.

## References

- [1] I. BOCK, *On nonstationary von Kármán equations*, Z. Angew. Math. Mech. 76 (1996), pp. 559–571.
- [2] I. BOCK, *On the semidiscretization and linearization of pseudoparabolic von Kármán system for viscoelastic plates*, Math. Meth. Appl. Sci. 29 (2006), pp. 557–573.
- [3] P.G. CIARLET and P. RABIER: *Les équations de von Kármán*. Springer, Berlin 1980.
- [4] C. ECK, J. JARUŠEK and M. KRBEČ: *Unilateral Contact Problems in Mechanics. Variational Methods and Existence Theorems*. Monographs & Textbooks in Pure & Appl. Math. No. 270.
- [5] J. JARUŠEK, J. MÁLEK, J. NEČAS and V. ŠVERÁK: *Variational inequality for a viscous drum vibrating in the presence of an obstacle*. *Rend. Mat.*, Ser. VII, **12** (1992), 943–958.

- [6] H. KOCH and A. STACHEL: Global existence of classical solutions to the dynamical von Kármán equations. *Math. Methods in Applied Sciences* 16 (1993), 581–586.
- [7] J. NEČAS and I. HLAVÁČEK: Mathematical Theory of Elastic and Elasto-Plastic Bodies: An Introduction. Elsevier, Amsterdam-Oxford-New York 1981.
- [8] J.E. MUÑOZ RIVERA and G.P. MENZALA: Decay rates of solutions to a von Kármán system for viscoelastic plates with memory. *Quarterly Appl. Math.* 57, 1 (1) (1999), 181–200.

### Current address

#### **Mária Kečkemétyová, RNDr. PhD.**

Department of Mathematics, Faculty of Electrical Engineering and Information Technology,  
Slovak University of Technology, Ilkovičova 3, 812 19 Bratislava 1, Slovak Republic;

tel.number +421260291566,

e-mail: maria.keckemetyova@stuba.sk

#### **Igor Bock, Professor, PhD.**

Department of Mathematics, Faculty of Electrical Engineering and Information Technology,  
Slovak University of Technology, Ilkovičova 3, 812 19 Bratislava 1, Slovak Republic;

tel.number +421260291204,

e-mail: igor.bock@stuba.sk



## STRICT FIXED POINT PRINCIPLES AND APPLICATIONS TO MATHEMATICAL ECONOMICS

MUREȘAN Anton S., (RO)

**Abstract.** In this paper we give some new results on strict fixed points for multivalued operators. We prove that there exists at least a strict fixed point, and we give conditions which assure that if the strict fixed point set of a multivalued operator is nonempty then the fixed point set and the strict fixed point set are equal. Some applications to mathematical economics are given too.

**Key words and phrases.** multivalued operators, fixed points, strict fixed points, abstract economy

*Mathematics Subject Classification.* 47H10, 54H25, 90B15, 91B52

### 1 Introduction

In the theory of fixed points for singlevalued or multivalued operators have been given several results.

The theory of strict fixed points of multivalued operators has been less investigated. After knowledge of author, a synthesis book having as an exclusive subject the strict fixed points does not exist.

However, some important results on strict fixed points were obtained by many authors like: I.A. Rus [20],[22]-[24], A. Petrușel [17], A. Muntean [10], A.S. Mureșan [11], A. Sîntămărian [25], A. Avram [2], A. Ahmad & M. Imbad [1], T.L. Hicks [5], K. Iseki [6], T. Kubiak [7], N. Negoescu [15], D.H. Tan & D.T. Nhan [27].

Some mathematical economics applications are given by H.W. Corley [4], A. Muntean [10] and J.X. Zhou [28].

In Section 2 we give some basic notions and needed results on strict fixed points in metric spaces.

Section 3 is dedicated to the conditions which assure that strict fixed points set of a multivalued operator is a nonempty set.

In Section 4 we give some conditions which assure that if strict fixed points set is nonempty then the fixed points set and strict fixed points set are equal.

Some applications to mathematical economics are given in Section 5.

The aim of this paper is to present some new results on strict fixed points for multivalued operators on metric spaces and some applications to mathematical economics.

## 2 Basic notions and needed results on fixed and strict fixed points

Let  $X$  be a nonempty set and  $T : X \multimap X$  a multivalued operator.

**Definition 2.1** An element  $x \in X$  is

- a **fixed point** of  $T$  iff  $x \in T(x)$ ;
- a **strict fixed point** of  $T$  iff  $T(x) = \{x\}$ ;
- a **univalent point** of  $T$  iff  $\text{card } T(x) = 1$ .

We denote by  $F_T$  the fixed points set of  $T$ , by  $(SF)_T$  the strict fixed points set of  $T$  and by  $U_T$  the univalent points set of  $T$ .

**Remark 2.2** A strict fixed point of  $T$  is a univalent point of  $T$ , therefore  $(SF)_T \subset U_T$ . The conversely assertion is not true.

**Definition 2.3** A subset  $A \subset X$  is

- an **invariant subset** of  $T$  iff  $T(A) \subset A$
- a **fixed subset** of  $T$  if  $T(A) = A$ .

We denote by  $I(T) := \{A \mid A \in P(X), T(A) \subset A\}$  the set of invariant subsets of  $T$  and by  $F(T) = \{A \mid A \in P(X), T(A) = A\}$  the set of fixed points subset of  $T$ .

We have (see [20]):

**Theorem 2.4** Let  $T : X \multimap P(X)$  a multivalued operator. Then

- a)  $T^n(X)$  is an invariant subset of  $T$ ,  $n \in \mathbb{N}$ ,  $T^n(X) \in I(T)$ .
- b)  $(SF)_T$  is a fixed subset of  $T$ ,  $(SF)_T \in F(T)$ .

**Remark 2.5** For any multivalued operator  $T$ , we have  $F_T \subset T(F_T)$ .

**Remark 2.6** Generally,  $F_T$  is not an invariant subset of  $T$ .



We give some conditions in which  $F_T$  is an invariant subset of  $T$ , therefore we have  $T(F_T) = F_T$ , or  $F_T \in F(T)$ .

**Theorem 2.7** *Let  $X$  be a nonempty set and  $T : X \multimap X$  a multivalued operator. If at least one from the following conditions are satisfied:*

- a)  $F_T \subset U_T$ ;
  - b) for any  $A \subset X$  there exists  $n \in \mathbb{N}$  such that  $T^n(A) \subset F_T$ ;
  - c) for any  $x \in X, T(x) \subset F_T$ ,
- then  $T(F_T) = F_T$ .

**Proof.** It is enough to prove that  $F_T \in I(T)$ , that is,  $F_T$  is an invariant subset of  $T$ .

a) If  $x \in F_T$  then  $x \in T(x)$ . But, because  $F_T \subset U_T$  it results that  $x \in U_T$ , therefore  $\text{card } T(x) = 1$ . Thus  $T(x) = \{x\}$  for any  $x \in F_T$ . It means that  $F_T = (SF)_T$ . Then, from Lemma 2.1. [11] we have  $T(F_T) = T((SF)_T) = (SF)_T = F_T$ .

b) For  $x \in T(F_T)$  there exists  $y \in F_T$  such that  $x \in T(y)$ . Because  $T(y) \subset X$  there exists  $n \in \mathbb{N}$  such that  $T^n(T(y)) = T^{n+1}(y) \subset F_T$ .

Because  $y \in F_T$  we have  $y \in T(y)$  therefore  $x \in T(y) \subset T(T(y)) = T^2(y)$ . Analogously, we obtain that  $x \in T^{n+1}(y)$ . Thus  $x \in F_T$ , hence  $T(F_T) \subset F_T$ .

c) For  $x \in T(F_T)$  there exists  $y \in F_T$  such that  $x \in T(y)$ . But  $T(y) \subset F_T$ , therefore  $x \in F_T$ , hence  $T(F_T) \subset F_T$ .

The theorem is proved.

Let  $(X, d)$  be a metric space. We denote:

$$P(X) := \{Y | \emptyset \neq Y \subset X\}$$

$$P_p(X) := \{Y | Y \in P(X), Y \text{ has the property } p\},$$

where  $p$  could be:  $b$  = bounded,  $cl$  = closed,  $cp$  = compact, ... .

We consider, in what follows, the following functionals:

$$D : P(X) \times P(X) \rightarrow \mathbb{R}_+, D(Y, Z) := \inf\{d(y, z) | y \in Y, z \in Z\}$$

$$\delta : P(X) \times P(X) \rightarrow \mathbb{R}_+ \cup \{+\infty\}, \delta(Y, Z) := \sup\{d(y, z) | y \in Y, z \in Z\}$$

$$\rho : P(X) \times P(X) \rightarrow \mathbb{R}_+ \cup \{+\infty\}, \rho(Y, Z) := \sup\{D(y, Z) | y \in Y\}$$

$$H : P(X) \times P(X) \rightarrow \mathbb{R}_+ \cup \{+\infty\}, H(Y, Z) := \sup\{\rho(Y, Z), \rho(Z, Y)\}$$

For the basic properties of these functionals see [22] and [10].

**Theorem 2.8** Let  $(X, d)$  be a complete metric space,  $T : X \rightarrow P_b(X)$  and  $\varphi : \mathbb{R}_+^5 \rightarrow \mathbb{R}_+$ . We suppose that:

- i)  $r, s \in \mathbb{R}_+^5$ ,  $r \leq s$  implies  $\varphi(r) \leq \varphi(s)$ ;
- ii)  $\varphi(u, u, u, u, u) < u$ , for all  $u \in \mathbb{R}, u > 0$ ;
- iii) for all  $x, y \in X$ ,

$$\delta(T(x), T(y)) \leq \varphi(d(x, y), \delta(x, T(x)), \delta(y, T(y)), \delta(x, T(y)), \delta(y, T(x))).$$

Then  $F_T = (SF)_T$ .

**Proof.** We have  $(SF)_T \subset F_T$ .

If  $F_T = \emptyset$  then  $(SF)_T = \emptyset$  and so  $F_T = (SF)_T$ .

Let now  $F_T \neq \emptyset$  and let  $x \in F_T$ . With this  $x$ , and  $y = x$  in iii), we have

$$\delta(x, x) = \delta(T(x), T(x)) \leq \varphi(d(x, x), \delta(x, T(x)), \delta(x, T(x)), \delta(x, T(x)), \delta(x, T(x))) \leq$$

$$\leq \varphi(\delta(x, T(x)), \delta(x, T(x)), \delta(x, T(x)), \delta(x, T(x)), \delta(x, T(x))) < \delta(x, T(x)),$$

if  $\delta(x, T(x)) > 0$ , that is impossible.

Therefore  $\delta(x, T(x)) = 0$ , that is  $T(x) = \{x\}$ . Thus  $x \in (SF)_T$  and  $F_T \subset (SF)_T$ . We obtain that  $F_T = (SF)_T$ .

The theorem is proved.

**Theorem 2.9** ([23]) Let  $(X, d)$  be a complete metric space,  $T : X \rightarrow P_b(X)$  and  $\varphi : \mathbb{R}_+^5 \rightarrow \mathbb{R}_+$ . We suppose that:

- i)  $r, s \in \mathbb{R}_+^5$ ,  $r \leq s$  implies  $\varphi(r) \leq \varphi(s)$ ;
- ii) there exists  $p > 1$  such that  $\varphi(u, pu, pu, u, u) < u$ , for all  $u \in \mathbb{R}, u > 0$ ;
- iii)  $u - \varphi(u, pu, pu, u, u) \rightarrow +\infty$  as  $u \rightarrow +\infty$ ;
- iv)  $\varphi$  is continuous;
- v) for all  $x, y \in X$ ,

$$\delta(T(x), T(y)) \leq \varphi(d(x, y), \delta(x, T(x)), \delta(y, T(y)), D(x, T(y)), D(y, T(x))).$$

Then  $F_T = (SF)_T = \{x^*\}$ .

**Proof.** Let  $p > 1$ . By Lemma 8.1.3 in [22] there exists a selection  $t$  of  $T$  such that  $\delta(x, T(x)) \leq p d(x, t(x))$  for all  $x \in X$ . From condition v) it follows that

$$d(t(x), t(y)) \leq \varphi(d(x, y), \delta(x, t(x)), \delta(y, t(y)), d(x, t(y)), d(y, t(x))).$$

This means that the selection operator  $t$  is a  $\varphi$ -contraction, such that there exists a unique fixed point  $x^*$  for  $t$ . So  $F_T \neq \emptyset$ . Let  $x \in F_T$ . If we take  $y = x$  in v), it results that

$$\delta(T(x)) = \delta(T(x), T(x)) \leq \varphi(0, \delta(x, T(x)), \delta(x, T(x)), 0, 0) \leq$$

$$\leq \varphi(\delta(T(x)), p\delta(T(x)), p\delta(T(x)), \delta(T(x)), \delta(T(x))) < \delta(T(x))$$

if  $\delta(T(x)) > 0$ , that is impossible. So, we have  $\delta(T(x)) = 0$ , that is  $T(x) = \{x\}$ , hence  $F_T = (SF)_T \neq \emptyset$ .

The uniqueness of the strict fixed point follows from v).

The theorem is proved.

### 3 Strict fixed point set of a multivalued operator

Let  $(X, d)$  be a complete metric space and  $T : X \rightarrow P_b(X)$  a multivalued operator.

**Theorem 3.1** *If the following conditions are satisfied*

- i) for all  $x \in X$ ,  $x \in T(x)$ ;
- ii) there exists a comparison function  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  and a Picard sequence  $(x_n)_{n \in \mathbb{N}}$ ,  $x_{n+1} \in T(x_n)$ ,  $n \in \mathbb{N}$ , such that

$$\delta(T(x_{n+1})) \leq \varphi(\delta(T(x_n))), \quad n \in \mathbb{N},$$

then there exists  $x^* \in X$  such that  $x_n \rightarrow x^*$  as  $n \rightarrow +\infty$ , and  $x^* \in (SF)_T \neq \emptyset$ .

**Proof.** From condition ii) it follows that

$$\delta(T(x_n)) \leq \varphi^n(\delta(T(x_0))) \rightarrow 0 \quad \text{as } n \rightarrow +\infty.$$

This implies that  $(x_n)_{n \in \mathbb{N}}$  is a Cauchy sequence. So,  $(x_n)_{n \in \mathbb{N}}$  converges. Let  $x^* \in X$  be the limit of this sequence. From i) we have  $x^* \in T(x^*) \neq \emptyset$ , and from  $\delta(T(x^*)) = 0$  it results that  $x^* \in (SF)_T \neq \emptyset$ .

The theorem is proved.

**Remark 3.2** *If we take  $\varphi(t) = at$ ,  $0 \leq a < 1$ , then we obtain a result given by H.W. Corley in [4].*

**Theorem 3.3** *Let  $(X, d)$  be a complete metric space,  $T : X \rightarrow P_{cl}(X)$  a continuous operator and  $\varphi : \mathbb{R}_+^3 \rightarrow \mathbb{R}_+$  a function. We suppose that the following conditions are satisfied:*

- i)  $r, s \in \mathbb{R}_+^3$ ,  $r \leq s$  implies  $\varphi(r) \leq \varphi(s)$ ;
- ii)  $\varphi(u, u, u) \leq u$ , for all  $u \in \mathbb{R}$ ,  $u > 0$ ;
- iii) for all  $x, y \in X$ ,  $x \neq y$ ,

$$\delta^2(T(x), T^2(y)) < \varphi(d^2(x, y), H(x, T(x)) \cdot H(y, T^2(y)), D(x, T^2(y)) \cdot D(y, T(x))).$$

Then there exists  $x^* \in X$  such that  $x^* \in (SF)_T \neq \emptyset$  or  $x^* \in (SF)_{T^2} \neq \emptyset$ .

**Proof.** Because  $T$  is a continuous operator we can define the continuous functional  $f : X \rightarrow \mathbb{R}_+$ ,  $f(x) = H(x, T(x))$ . It follows that  $f$  takes its minimum value on  $X$ , hence there exists  $x_0 \in X$  such that  $f(x_0) = \inf\{f(x) | x \in X\}$ . We prove that  $x_0$  is a fixed point of  $T$  or some  $x_1 \in T(x_0)$  is a fixed point of  $T^2$ . We choose

$$x_1 \in T(x_0) \text{ such that } d(x_0, x_1) = H(x_0, T(x_0)),$$

$$x_2 \in T(x_1) \text{ such that } d(x_1, x_2) = H(x_1, T^2(x_1)),$$

$$x_3 \in T(x_2) \text{ such that } d(x_2, x_3) = H(x_2, T(x_2)).$$

We shall prove that  $H(x_0, T(x_0)) = 0$  or  $H(x_1, T^2(x_1)) = 0$ , that is  $T(x_0) = \{x_0\}$  or  $T^2(x_1) = \{x_1\}$ .

Suppose that  $H(x_0, T(x_0)) > 0$  and  $H(x_1, T^2(x_1)) > 0$ .

By using iii), we obtain

$$\begin{aligned} d^2(x_1, x_2) &\leq H^2(T(x_0), T^2(x_1)) < \\ &< \varphi(d^2(x_0, x_1), H(x_0, T(x_0)) \cdot H(x_1, T^2(x_1)), D(x_0, T^2(x_1)) \cdot D(x_1, T(x_0))) = \\ &= \varphi(d^2(x_0, x_1), d(x_0, x_1) \cdot d(x_1, x_2), 0) \leq \max\{d^2(x_0, x_1), d(x_0, x_1) \cdot d(x_1, x_2), 0\}. \end{aligned}$$

If  $d^2(x_1, x_2) < d^2(x_0, x_1)$  it follows that  $d(x_1, x_2) < d(x_0, x_1)$ .

If  $d^2(x_1, x_2) < d(x_0, x_1) \cdot d(x_1, x_2)$ , how  $d(x_1, x_2) = H(x_1, T^2(x_1)) > 0$ , it follows that  $d(x_1, x_2) < d(x_0, x_1)$ , the same inequality as before.

Then, we have

$$\begin{aligned} d^2(x_2, x_3) &\leq H^2(T(x_2), T^2(x_1)) < \\ &< \varphi(d^2(x_1, x_2), H(x_2, T(x_2)) \cdot H(x_1, T^2(x_1)), D(x_2, T^2(x_2)) \cdot D(x_1, T(x_2))) = \\ &= \varphi(d^2(x_1, x_2), d(x_2, x_3) \cdot d(x_1, x_2), 0) \leq \max\{d^2(x_1, x_2), d(x_2, x_3) \cdot d(x_1, x_2), 0\}. \end{aligned}$$

Similarly, it follows that

$$d^2(x_2, x_3) < d^2(x_1, x_2), \text{ such that } d(x_2, x_3) < d(x_1, x_2), \text{ or } d^2(x_2, x_3) < d(x_2, x_3) \cdot d(x_1, x_2).$$

In the second situation, if  $d(x_2, x_3) = 0$  we obtain a contradiction, such that we must have the same inequality  $d(x_2, x_3) < d(x_1, x_2)$ .

Therefore, we deduce successively that

$$H(x_2, T(x_2)) = d(x_2, x_3) < d(x_1, x_2) < d(x_0, x_1) = H(x_0, T(x_0)) = f(x_0)$$

which contradicts the minimality of  $f(x_0)$ .

So, we must have  $H(x_0, T(x_0)) = 0$ , that is  $T(x_0) = \{x_0\}$ , or  $H(x_1, T^2(x_1)) = 0$ , that is  $T^2(x_1) = \{x_1\}$ .

The proof is complete.

**Remark 3.4** A similar result can be obtained in the case when  $\varphi : \mathbb{R}_+^3 \rightarrow \mathbb{R}_+$ ,  $\varphi(t_1, t_2, t_3) = \max\{t_1, t_2, t_3\}$ . This result was given by N. Negoescu in [15].

#### 4 The case when the fixed set and the strict fixed point set are equal

**Theorem 4.1** Let  $(X, d)$  be a complete metric space and  $T : X \rightarrow P(X)$  a  $(\delta, \varphi)$ -contraction. Then there exists  $x^* \in X$  such that

$$F_T = (SF)_T = \{x^*\}.$$

**Proof.** The operator  $T$  is a  $(\delta, \varphi)$ -contraction. This means that there exists a comparison function  $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  such that

$$\delta(T(Y)) \leq \varphi(\delta(Y)), \quad \text{for all } Y \in I(T).$$

Let be the following sequence of sets, defined by

$$X_1 := \overline{T(X)}, \dots, X_{n+1} := \overline{T(X_n)}, \quad n \in \mathbb{N}^*, \quad X_0 := X.$$

This sequence has the following properties:

- a)  $X \supset X_1 \supset \dots \supset X_n \supset \dots$
- b)  $X_n \in P_{b,cl}(X)$  and  $X_n \in I(T)$
- c)  $\delta(X_n) \leq \varphi(\delta(X)) \rightarrow 0$ , as  $n \rightarrow +\infty$ .

We denote by  $X_\infty := \bigcap_{n \in \mathbb{N}} X_n$ .

From a), b), c) we have that  $X_\infty \in I(T)$  and  $\delta(X_\infty) = 0$ .

Because  $T(X_\infty) \in P(X)$  it follows that  $T(X_\infty) \neq \emptyset$ .

Hence there exists  $x^* \in X$  such that  $X_\infty = \{x^*\}$  and  $x^* \in (SF)_T$ .

But, on the other hand,  $F_T \subset \bigcap_{n \in \mathbb{N}} X_n = X_\infty$ .

These imply that  $F_T = (SF)_T = \{x^*\}$ .

The theorem is proved.

**Remark 4.2** If  $(X, d)$  is a bounded complete metric space then the Theorem 4.1 implies the Theorem 3.1.

**Theorem 4.3** Let  $(X, d)$  be a complete metric space,  $T : X \rightarrow P_{b,cl}(X)$  an operator and  $\varphi : \mathbb{R}_+^3 \rightarrow \mathbb{R}_+$  a continuous function. We suppose that the following conditions hold:

- i)  $r, s \in \mathbb{R}_+$ ,  $r \leq s$  implies  $\varphi(r) \leq \varphi(s)$ ;
- ii)  $\varphi(u, \alpha u, \beta u) < u$ , for all  $u \in \mathbb{R}_+$ , where  $\alpha, \beta \in \{0, 1, 2\}$  and  $\alpha + \beta = 2$ ;
- iii) for all  $x, y \in X$

$$\delta(T^2(x), T^2(y)) < \varphi(D(T(x), T(y)), D(T(y), T^2(x)), D(T(x), T^2(y)));$$

- iv)  $T$  is u.s.c. on  $X$ .

Then there exists  $x^* \in X$  such that

$$F_T = (SF)_T = \{x^*\}.$$

**Proof.** Let  $x_0 \in X$  be. We consider the sequence  $(x_n)_{n \in \mathbb{N}}$  obtained as follows:

$$x_{n+1} \in T(x_n) \text{ such that } d(x_n, x_{n+1}) = \delta(x_n, T(x_n)) := b_n.$$

This real sequence  $(b_n)_{n \in \mathbb{N}}$  is decreasing. Indeed, for  $n \geq 2$ , we have

$$\begin{aligned} b_n &= \delta(x_n, T(x_n)) \leq \delta(T^2(x_{n-1}), T^2(x_{n-2})) \leq \\ &\leq \varphi(D(T(x_{n-1}), T(x_{n-2})), D(T(x_{n-2}), T^2(x_{n-1})), D(T(x_{n-1}), T^2(x_{n-2}))) \leq \\ &\leq \varphi(d(x_n, x_{n-1}), d(x_n, x_{n+1}), d(x_n, x_n)) \leq \varphi(b_{n-1}, b_{n-1} + b_n, 0). \end{aligned}$$

If  $b_n > b_{n-1}$ , then we get a contradiction,

$$b_n \leq \varphi(b_n, 2b_n, 0) < b_n.$$

So,  $b_n \leq b_{n-1}$ ,  $n \geq 2$ .

Let us prove that  $b_n \rightarrow 0$  as  $n \rightarrow +\infty$ .

If  $\psi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  is defined by  $\psi(t) := \varphi(t, \alpha t, \beta t)$ , where  $\alpha, \beta \in \{0, 1, 2\}$  and  $\alpha + \beta = 2$ , then the function  $\psi$  is increasing and satisfies the condition  $\psi(t) < t$ , for all  $t > 0$ .

Inductively, we can prove that the relationships

$$b_n \leq \psi^{n-1}(b_1), \quad n > 1$$

hold.

Thus  $b_n \rightarrow 0$  as  $n \rightarrow +\infty$ . Easily, it results that  $(x_n)_{n \in \mathbb{N}}$  is a Cauchy sequence in the complete metric space  $(X, d)$ . Let  $x^*$  be the limit of this sequence.

The operator  $T$  is u.s.c. with closed and bounded values, such that  $T$  is a closed operator on  $X$ .

We have  $\delta(x^*, T(x^*)) = 0$  and  $\delta(x^*, T^2(x^*)) = 0$ . So  $T(x^*) = \{x^*\}$ ,  $T^2(x^*) = \{x^*\}$ , that is  $T(x^*) = T^2(x^*) = \{x^*\}$ , or  $(SF)_T \cap (SF)_{T^2} \neq \emptyset$ .

The uniqueness of common strict fixed point of  $T$  and  $T^2$  results easily.

We suppose that there exists  $y^* \in (SF)_T \cap (SF)_{T^2}$ , such that  $x^* \neq y^*$ . Then

$$d(x^*, y^*) = \delta(T^2(x^*), T^2(y^*)) \leq \varphi(d(x^*, y^*), d(x^*, y^*), d(x^*, y^*)) < d(x^*, y^*)$$

which is a contradiction. It follows that  $(SF)_T \cap (SF)_{T^2} = \{x^*\}$ .

Let us prove that  $F_T \cap F_{T^2} = (SF)_T \cap (SF)_{T^2}$ .

Suppose that there exists  $y^* \in F_T \cap F_{T^2}$  such that  $x^* \neq y^*$ .

Then, we have  $y^* \in F_{T^2}$  and therefore

$$\begin{aligned} d(x^*, y^*) &\leq \delta(T^2(x^*), T^2(y^*)) \leq \\ &\leq \varphi(D(T(x^*), T(y^*)), D(T(y^*), T^2(x^*)), D(T(x^*), T^2(y^*))) = \\ &= \varphi(D(x^*, T(y^*)), D(x^*, T(y^*)), d(x^*, y^*)) \leq \\ &\leq \varphi(d(x^*, y^*), d(x^*, y^*), d(x^*, y^*)) < d(x^*, y^*), \end{aligned}$$

that is a contradiction. This means that  $y^* = x^* \in (SF)_T \cap (SF)_{T^2}$ , and hence  $F_T \cap F_{T^2} = (SF)_T \cap (SF)_{T^2} = \{x^*\} = F_T = (SF)_T$ .

The theorem is proved.

## 5 Some applications to mathematical economics

The consumer's problem can be modelled in terms of strict fixed points of a multivalued operator.

We will adopt the classical line of the theory of J. von Neumann, K. Arrow and G. Debreu for abstract economies.

**Definition 5.1** Let  $(E)$  be an abstract economy with  $m$  consumers and  $n$  commodities. Then, by definition:

- i) the **consumer's set** is  $I = \{1, 2, \dots, m\}$ ,
- ii) the **commodity space** is  $\mathbb{R}^n$ ,
- iii) the **consumption set** (the possible consumption set or the **choice set**) for each consumer  $i \in I$ , is a set  $Y_i \subset \mathbb{R}^n$ ,
- iv) a **possible consumption vector** for the consumer  $i$  (the  $i^{\text{th}}$  **consumer's demand**) is an  $n$ -dimensional vector  $x_i \in Y_i$ ,
- v) the **prices' simplex** (the set of admissible prices) is

$$\sigma^n := \{p | p \in \mathbb{R}^n, \sum_{k=1}^n p_k = 1, p_k > 0\}.$$

**Definition 5.2** For each consumer  $i \in I$  we define:

- i) the **budget set**  $B_i(p) := \{x_i | x_i \in Y_i, p \cdot x_i \leq 1\}$ ,
- ii) the **choosing multivalued operator** (the **preference operator**)

$$U_i : Y_i \multimap Y_i, \quad U_i(x_i) := \{y_i | y_i \in Y_i, y_i \succeq x_i\}$$

- iii) the **optimal preference** is a consumption vector  $x_i^* \in Y_i$  satisfying the condition  $U_i(x_i^*) = \{x_i^*\}$ .

**Definition 5.3** The **consumer's problem** consists in the choice of a consumption vector  $x_i^* \in B_i(p)$  such that to have  $x_i^* \succeq x_i$ , for all  $x_i \in B_i(p)$ . So, a consumption vector  $x_i^*$  is chosen as the "optimal" variant by the consumer, if  $U_i(x_i^*) = \{x_i^*\}$ .

**Definition 5.4** Let  $Y$  be a topological vector space. A set  $C \in P(Y)$  is called **cone** iff for all  $x \in C$ , and  $0 < \lambda \in \mathbb{R}$  it results  $\lambda \cdot x \in C$ .

**Definition 5.5** A cone  $C \subset Y$  is

- i) a **convex cone** if it is a convex set;
- ii) a **pointed cone** if  $C \cap (-C) = \emptyset$ ;
- iii) an **acute cone** if  $\overline{C}$  is a pointed cone.

The following result holds

**Theorem 5.6** ([10]) Let  $Y_i \in P(\mathbb{R}^n)$ ,  $C \subset \mathbb{R}^n$  be an acute convex cone and  $U_i : Y_i \rightarrow P(Y_i)$  the preference operator. We suppose that:

- i)  $Y_i$  is a  $C$ -semicompact set;
- ii) for all  $x_i \in Y_i$ , the upper contour set  $U_i(x_i)$  is a  $C$ -semicompact set.

Then the consumer's problem has at least one solution, that is, there exists  $x_i^* \in B_i(p)$  such that  $U_i(x_i^*) = \{x_i^*\}$ .

**Definition 5.7** An **abstract economy** (or a **generalized game**) is defined as a family of ordered triples  $\Gamma = (Y_i, F_i, U_i)_{i \in I}$  where  $Y_i$  is a **choice set** (a nonempty topological vector space),  $F_i : Y \multimap Y_i$ ,  $Y := \prod_{i \in I} Y_i$ , are **constraint multivalued operator** and  $U_i$  is the **preference operator**. An **equilibrium choice** for  $\Gamma$  (of Schafer-Sonnenschein type) is a point  $x^* \in Y$  such that for each  $i \in I$ ,  $x_i^* \in \text{cl } F_i(x^*)$  and  $U_i(x^*) \cap F_i(x^*) = \emptyset$ .

**Theorem 5.8** Let  $\Gamma = (Y_i, F_i, U_i)_{i \in I}$  be an abstract economy such that for each  $i \in I$ ,

- 1)  $Y_i$  is a nonempty compact convex subset of a metrisable locally convex Hausdorff topological vector space,
- 2) for each  $x \in Y$ ,  $F_i(x)$  is a nonempty convex subset of  $Y_i$ ,
- 3) the multivalued operator  $\text{cl } F_i : Y \multimap Y_i$  is continuous,
- 4) the multivalued operator  $U_i$  is  $\Theta$ -majorised.

Then the abstract economy  $\Gamma$  has an equilibrium choice  $x^* \in Y$ , that is, for each  $i \in I$ ,  $x_i^* \in \text{cl } F_i(x^*)$  and  $U_i(x^*) \cap F_i(x^*) = \emptyset$ .

**Proof.** Let  $i \in I$  be fixed. Since  $U_i$  is  $\Theta$ -majorised, for each  $x \in Y$ , there exists a multivalued operator  $\phi_x : Y \multimap Y_i$  and an open neighborhood  $N_x$  of  $x$  in  $Y$  such that  $U_i(z) \subset \phi_x(z)$  and  $z_i \notin \text{cl co } \phi_x(z)$  for each  $z \in N_x$ , and  $\phi_x|_{N_x}$  has an open graph in  $N_x \times Y_i$ . By compactness of  $Y$ , the family  $\{N_x | x \in Y\}$  is an open cover of  $Y$  which contains a finite subcover  $\{N_{x_j} | j \in J\}$ , where  $J$  is a finite set,  $J \subset \mathbb{N}$ . For each  $j \in J$ , we now define  $\phi_j : Y \multimap Y_i$  by  $\phi_j(z) = \phi_{x_j}(z)$  if  $z \in N_{x_j}$ , respectively  $\phi_j(z) = Y_i$  if  $z \notin N_{x_j}$ , and next we define  $\Phi_i : Y \multimap Y_i$  by  $\Phi_i(z) = \bigcap_{j \in J} \phi_j(z)$ . For each  $z \in Y$ , there exists  $k \in J$  such that  $z \in N_{x_k}$  and so that  $z_i \notin \text{cl co } \phi_{x_k}(z)$ . Thus  $z_i \notin \text{cl co } \Phi_i(z)$ .



We now show that the graph of  $\Phi_i$  is open in  $Y \times Y_i$ . For each  $(z, x) \in \text{graph of } \Phi_i$ , since  $Y = \bigcup_{j \in J} N_{x_j}$ , there exists  $\{i_1, \dots, i_k\} \subset J$  such that  $z \in N_{x_{i_1}} \cap \dots \cap N_{x_{i_k}}$ . Then we can find an open neighborhood  $N$  of  $z$  in  $Y$  such that  $N \subset N_{x_{i_1}} \cap \dots \cap N_{x_{i_k}}$ . Since  $\phi_{x_{i_1}}(z) \cap \dots \cap \phi_{x_{i_k}}(z)$  is an open subset of  $Y_i$  containing  $x$ , there exists an open neighborhood  $V$  of  $x$  in  $Y_i$  such that  $x \in V \subset \phi_{x_{i_1}}(z) \cap \dots \cap \phi_{x_{i_k}}(z)$ . Therefore we have an open neighborhood  $N \times V$  of  $(z, x)$  such that  $N \times V \subset \text{graph of } \Phi_i$ , so that the graph of  $\Phi_i$  is open in  $Y \times Y_i$ . And it is clear that  $U_i(z) \subset \Phi_i(z)$  for each  $z \in Y$ .

Next, since  $Y \times Y_i$  is compact and metrisable, so it is perfectly normal. Since the graph of  $\Phi_i$  is open in  $Y \times Y_i$ , by a result of Dugundji, there exists a continuous function  $C_i : Y \times Y_i \rightarrow [0, 1]$  such that  $C_i(x, y) = 0$  for all  $(x, y) \notin \text{graph of } \Phi_i$  and  $C_i(x, y) \neq 0$  for all  $(x, y) \in \text{graph of } \Phi_i$ . For each  $i \in I$  we define the multivalued operator  $G_i : Y \multimap Y_i$  by

$$G_i(x) = \{y | y \in cl F_i(x), C_i(x, y) = \max_{z \in cl F_i(x)} C_i(x, z)\}.$$

Then by a result of Aubin and Ekeland,  $G_i$  is upper semicontinuous and for each  $x \in Y$ ,  $G_i(x)$  is nonempty closed subset of  $Y_i$ . Then the multivalued operator  $G : Y \multimap Y$  defined by  $G(x) = \prod_{i \in I} G_i(x)$  is also upper semicontinuous, by a result of Fan, and  $G(x)$  is a nonempty compact subset of  $Y$ , for each  $x \in Y$ . Therefore, by the Corollary 3.1. [13], there exists a point  $x^* \in Y$  such that  $x^* \in cl co G(x^*)$ , that is,  $x^* \in cl co G(x^*) \subset \prod_{i \in I} cl co G_i(x^*)$ . Since  $G_i(x^*) \subset cl F_i(x^*)$  and  $F_i(x^*)$  is convex,  $cl co G_i(x^*) \subset cl F_i(x^*)$ . Therefore  $x_i^* \in cl F_i(x^*)$  for each  $i \in I$ . It remains to show that  $U_i(x^*) \cap F_i(x^*) = \emptyset$ . If  $z_i \in U_i(x^*) \cap F_i(x^*) \neq \emptyset$ , then  $C_i(x^*, z_i) > 0$  so that  $C_i(x^*, z'_i) > 0$  for all  $z'_i \in F_i(x^*)$ . This implies that  $F_i(x^*) \subset \Phi_i(x^*)$ , which implies  $x_i^* \in cl co F_i(x^*) \subset cl co \Phi_i(x^*)$ ; this is a contradiction. So, the theorem is proved.

## References

- [1] AHMAD, A., IMBAD, M., *Some common fixed point theorems for mappings and multivalued mappings*, J. Math. Anal. Appl., **218**(1998), 546-560
- [2] AVRAM, M., *Points fixed communs pour les applications multivoques dans les espaces metrique*, Mathematica, **17**(1975), 2, 153-156
- [3] BORDER, K.C., *Fixed point theorems with applications to economics and game theory*, Cambridge University Press, 1985
- [4] CORLEY, H.W., *Some hybrid fixed point theorem related to optimization*, J. Math. Anal. Appl., **120**(1986), 528-532
- [5] DEBREU, G., *Economies with a finite set of equilibria*, Econometrica, **38**(1970), 387-392
- [6] HICKS, T.L., *Setvalued mappings on metric spaces*, Indian J. Pure Appl. Math., **22**(1991), 269-271
- [7] ISEKI, K., *Multivalued contraction mappings in complete metric spaces*, Math. Sem. Notes, **2**(1974), 45-49
- [8] KUBIAK, T., *Fixed point theorems for contractive type multivalued mappings*, Math. Japonica, **30**(1985), 1, 89-101

- [9] MEHTA, G., SESSA, S., *Coincidence theorems and maximal elements in topological vector spaces*, Math. Japonica, **37**(1992), 839-845
- [10] MUNTEAN, A., *Fixed point principles and applications to mathematical economics*, Cluj University Press, Cluj-Napoca, 2002
- [11] MUREŞAN, A.S., *On some invariant problem of fixed points set for multivalued mappings*, Seminar on Fixed Point Theory, Cluj-Napoca, 1985, 37-42
- [12] MUREŞAN, A.S., *First order equilibria for an abstract economy*, Bull. Şt. Univ. Baia Mare, Ser. Math.-Inform., XIV (1998), 191-196
- [13] MUREŞAN, A.S., *Non-cooperative games*, Ed. Mediamira, Cluj-Napoca, 2003
- [14] NEGOESCU, N., *Observations sur des paires d'applications multivoques d'un certain type de contractivite*, Bul. Inst. Polit. Iaşi, **35**(1989), 21-25
- [15] NEGOESCU, N., *Common fixed points problem for pairs  $\varphi$ - contractive of functions (Romanian)*, Ed. Gh. Asachi, Iaşi, 1999
- [16] PETRUŞEL, A., *Multivalued operators in generalized metric spaces*, Seminar on Fixed Point Theory, Cluj-Napoca, 1994, 11-16
- [17] PETRUŞEL, A., *The multivalued operators' analysis*, Babeş-Bolyai University, Cluj-Napoca, 1996
- [18] PETRUŞEL, A., *Multivalued operators and continuous selection. The fixed point set*, PU. M.A., **9**(1998), 165-170
- [19] RUS, A.I., *Fixed point theorems for multivalued mappings in complete metric spaces*, Math. Japonica, **20**(1975), 21-24
- [20] RUS, A.I., *Fixed and strict fixed points for multivalued mappings*, Seminar on Fixed Point Theory, Cluj-Napoca, 1985, 77-82
- [21] RUS, A.I., *Technique of the fixed point structures for multivalued mappings*, Math. Japonica, **38**(1993), 289-296
- [22] RUS, A.I., *Generalized contractions and applications*, Cluj University Press, Cluj-Napoca, 2001
- [23] RUS, A.I., *Strict fixed point theory*, Fixed Point Theory, **4**(2003), 2, 177-183
- [24] RUS, A.I., PETRUŞEL, A., PETRUŞEL, G., *Fixed Point Theory 1950-2000. Romanian Contributions*, House of the Book of Science, Cluj-Napoca, 2002
- [25] SÎNTĂMĂRIAN, A., *Metrical strict fixed point theorems for multivalued mappings*, Seminar on Fixed Point Theory, Cluj-Napoca, 1997, 27-31
- [26] SÎNTĂMĂRIAN, A., *Common fixed point theorems for multivalued mappings*, Seminar on Fixed Point Theory, Cluj-Napoca, **1**(2001), 93-102
- [27] TAN, D.H., NHAN, D.T., *Common fixed points of two mappings of contractive type*, Acta Math. Vietnamica, **5**(1980), 150-160
- [28] ZHOU, J.X., *On the existence of equilibrium for abstract economies*, J. Math. Appl., **193**(1995), 839-858

#### **Current address**

**Anton S. MUREŞAN, prof.**

Department of Statistics, Forecasting and Mathematics,

Faculty of Economics and Business Administration  
Babeş-Bolyai University of Cluj-Napoca  
58-60 T. Mihali Street, Cluj-Napoca  
e-mail: [anton.muresan@econ.ubbcluj.ro](mailto:anton.muresan@econ.ubbcluj.ro)



# A FREDHOLM INTEGRAL EQUATION WITH LINEAR MODIFICATION OF THE ARGUMENT

MURESAN Viorica, (RO)

**Abstract.** In this paper we give existence and uniqueness results for the solution of a Fredholm integral equation with linear modification of the argument, in Banach space. We use Picard and weakly Picard operators' technique (see I.A.Rus [22]-[25]).

**Key words and phrases.** functional-integral equations, Picard operators, weakly Picard operators

*Mathematics Subject Classification.* 34K15, 34G20, 45N05, 47H10

## 1 Introduction

The theory of integral equations is an active field in mathematics. Many problems arising in natural and social sciences such as physics, mechanics, astronomy, chemistry, biology, economics, engineering lead to mathematical models described by functional integral equations. The theory of functional integral equations has developed very much. Many monographs appeared: Bellman and Cooke [2] (1963), Halanay [10] (1965), Elsgoltz and Norkin [5] (1971), Bernfeld and Lakshmikantham [3] (1974), Hale [9] (1977), Lakshmikantham [13] (1984), Azhelev, Maksimov and Rahmatulina [1] (1991), Hale and Verdyn Lunel [9] (1993), Guo and Lakshmikantham [7] (1996) such as a large number of papers. We quote here [11], [12], [18], [26], [27].

Let  $(X, || \cdot ||)$  be a Banach space. Consider the following Fredholm integral equation:

$$x(t) = g(t, x(t), x(\lambda t), x(0)) + \int_0^b K(t, s, x(s), x(\lambda s)) ds, \quad t \in [0, b], 0 < \lambda < 1, \quad (1.1)$$

where  $g \in C([0, b] \times X^3, X)$  and  $K \in C([0, b] \times [0, b] \times X^2, X)$ .

By using Picard and weakly Picard operators' technique we obtain existence and uniqueness results for the solution of the above equation and also comparison results.

## 2 Needed results from Picard operators theory

Let  $(X, d)$  be a metric space and  $A : X \rightarrow X$  an operator. We denote  
 $P(X) := \{Y \subseteq X | Y \neq \emptyset\}$ ;  $I(A) = \{Y \in P(X) | A(Y) \subseteq Y\}$ ;  
 $F_A := \{x \in X | A(x) = x\}$  - the fixed point set of  $A$ ;  
 $A^0 := 1_X, A^1 := A, A^{n+1} := A \circ A^n, n \in \mathbb{N}$ .

**Definition 2.1** (Rus [23])  $A$  is a **Picard operator** if there exists  $x^* \in X$  such that:

- (i)  $F_A = \{x^*\}$ ;
- (ii) the sequence  $(A^n(x_0))_{n \in \mathbb{N}}$  converges to  $x^*$ , for all  $x_0 \in X$ .

**Definition 2.2** (Rus [22], [24])  $A$  is a **weakly Picard operator** if the sequence  $(A^n(x_0))_{n \in \mathbb{N}}$  converges for all  $x_0 \in X$  and its limit (which may depend on  $x_0$ ) is a fixed point of  $A$ .

If  $A$  is a weakly Picard operator then we consider the operator  $A^\infty : X \rightarrow X$  defined by

$$A^\infty(x) = \lim_{n \rightarrow \infty} A^n(x), x \in X.$$

We have that  $A^\infty(X) = F_A$ .

**Remark 2.3** If  $A$  is a weakly Picard operator and  $F_A = \{x^*\}$ , then  $A$  is a Picard operator.

Let  $(X, d, \leq)$  be an ordered metric space and  $A : X \rightarrow X$  an operator.

We have:

**Lemma 2.4** (Rus [24]) If the operator  $A$  is increasing and  $A$  is a weakly Picard operator, then  $A^\infty$  is increasing.

**Lemma 2.5** (Comparison abstract lemma) (Rus [24]). Let  $A, B, C : X \rightarrow X$  three operators such that:

- (i)  $A, B, C$  are weakly Picard operators;
- (ii)  $B$  is an increasing operator;
- (iii)  $A \leq B \leq C$ .

Then  $x \leq y \leq z$  implies  $A(x) \leq B(y) \leq C(z)$ .

**Lemma 2.6** (Abstract Gronwall lemma) (Rus [24]) If the operator  $A$  is an increasing operator and  $A$  is a Picard operator with  $F_A = \{x^*\}$ , then

- (a)  $x \leq A(x)$  implies  $x \leq x^*$ ;
- (b)  $x \geq A(x)$  implies  $x \geq x^*$ .

### 3 Main results

Let  $(X, \|\cdot\|, \leq)$  be an ordered Banach space and  $(C([0, b], X), \|\cdot\|_C)$ , where  $\|x\|_C = \max_{t \in [0, b]} \|x(t)\|$ . Consider the equation (1.1) and we suppose that:

(i)  $g \in C([0, b] \times X^3, X)$ ,  $K \in C([0, b] \times [0, b] \times X^2, X)$ ;

(ii)  $K(0, s, u, v) = 0$ , for all  $s \in [0, b]$  and all  $u, v \in X$ .

Let  $S_g = \{\alpha \in X | g(0, \alpha, \alpha, \alpha) = \alpha\}$  be. We denote  $X_\alpha := \{x \in C([0, b], X) | x(0) = \alpha\}$ . It is clear that  $C([0, b], X) = \bigcup_{\alpha \in X} X_\alpha$  is a partition of  $C([0, b], X)$  and  $X_\alpha \in I(A)$  if and only if  $\alpha \in S_g$ .

Consider the operator  $A : C([0, b], X) \rightarrow C([0, b], X)$  defined by

$$A(x)(t) := g(t, x(t), x(\lambda t), x(0)) + \int_0^b K(t, s, x(s), x(\lambda s)) ds, \quad t \in [0, b], \quad 0 < \lambda < 1 \quad (3.1)$$

We have

**Theorem 3.1** (Theorem 3.1, [16]) *We suppose that the previous conditions (i) and (ii) are satisfied and*

(iii)  $S_g \neq \emptyset$

(iv) *there exists  $L_K > 0$ , such that*

$$\|K(t, s, u_1, v_1) - K(t, s, u_2, v_2)\| \leq L_K(\|u_1 - u_2\| + \|v_1 - v_2\|),$$

*for all  $t, s \in [0, b]$  and all  $u_i, v_i \in X$ ,  $i = 1, 2$ ;*

(v) *there exists  $0 < L_g < \frac{1}{2}$  such that*

$$\|g(t, u_1, u_2, \alpha) - g(t, v_1, v_2, \alpha)\| \leq L_g(\|u_1 - v_1\| + \|u_2 - v_2\|),$$

*for all  $t \in [0, b]$  and all  $u_i, v_i, \alpha \in X$ ,  $i = 1, 2$ .*

Then

$$A|_{\bigcup_{\alpha \in S_g} X_\alpha} : \bigcup_{\alpha \in S_g} X_\alpha \rightarrow \bigcup_{\alpha \in S_g} X_\alpha,$$

*is a weakly Picard operator.*

**Corollary 3.2** *Card  $F_A = \text{Card } S_g$ .*

**Theorem 3.3** *We suppose that (i) and (ii) are satisfied and*

(iii)  $S_g = \{\alpha^*\}$ ;

(iv) there exists  $L_K > 0$  such that

$$\|K(t, s, u_1, v_1) - K(t, s, u_2, v_2)\| \leq L_K(\|u_1 - u_2\| + \|v_1 - v_2\|),$$

for all  $t, s \in [0, b]$  and all  $u_i, v_i \in X$ ,  $i = 1, 2$ ;

(v) there exists  $L_g > 0$  such that

$$\|g(t, u_1, v_1, \alpha) - g(t, u_2, v_2, \alpha)\| \leq L_g(\|u_1 - u_2\| + \|v_1 - v_2\|),$$

for all  $t \in [0, b]$ , and all  $u_i, v_i, \alpha \in X$ ,  $i = 1, 2$ ;

(vi)  $L_g + L_K b < \frac{1}{2}$ .

Then the equation (1.1) has a unique solution  $x^*$  in  $C([0, b], X)$ .

**Proof.** Consider  $X_{\alpha^*} = \{x \in C([0, b], X) | x(0) = \alpha^*\}$  and  $A_* = A|_{X_{\alpha^*}}$ .

We have

$$\|A_{\alpha^*}(x)(t) - A_{\alpha^*}(z)(t)\| \leq \|g(t, x(t), x(\lambda t), \alpha^*) - g(t, z(t), z(\lambda t), \alpha^*)\| + \int_0^b \|K(t, s, x(s), x(\lambda s)) - K(t, s, z(s), z(\lambda s))\| ds \leq L_g(\|x(t) - z(t)\| + \|x(\lambda t) - z(\lambda t)\|) + 2L_K b \|x - z\|_C \leq 2(L_g + L_K b) \|x - z\|_C.$$

It follows that

$$\|A_{\alpha^*}(x) - A_{\alpha^*}(z)\|_C \leq 2(L_g + L_K b) \|x - z\|_C.$$

Because of (vi), the operator  $A_{\alpha^*}$  is a contraction. So,  $A^*$  is a Picard operator.

**Theorem 3.4** We suppose that all the conditions in Theorem 3.2. are satisfied and  $x^*$  is the unique solution of (1.1) and in addition we suppose that  $K(t, s, \cdot, \cdot, \cdot)$  and  $g(t, \cdot, \cdot, \alpha^*)$  are increasing, for all  $t, s \in [0, b]$ . In these conditions, if  $x$  is a subsolution of (1.1), then  $x \leq x^*$ .

**Proof.** We apply Lemma 2.3.

Consider the equations:

$$x(t) = g_i(t, x(t), x(\lambda t), x(0)) + \int_0^b K_i(t, s, x(s), x(\lambda s)) ds, \quad t \in [0, b], 0 < \lambda < 1, \quad (3.2)_i$$

where  $i = 1, 2, 3$ .

We have

**Theorem 3.5** We suppose that for (3.2)<sub>i</sub>,  $i = 1, 2, 3$  the corresponding conditions of Theorem 3.2 are satisfied and let  $x_i^*$  be,  $i = 1, 2, 3$  the corresponding solution for each equation. If, in addition,  $Sg_1 = Sg_2 = Sg_3 = \{\alpha^*\}$  and we suppose that  $K_2(t, s, \cdot, \cdot, \cdot)$  and  $g_2(t, \cdot, \cdot, \alpha^*)$  are increasing for all  $t, s \in [0, b]$  and  $K_1 \leq K_2 \leq K_3$ ,  $g_1 \leq g_2 \leq g_3$ , then  $x_1^* \leq x_2^* \leq x_3^*$ .

**Proof.** Consider the corresponding operators  $A_i$ ,  $i = 1, 2, 3$ , defined by

$$A_i(x)(t) := g_i(t, x(t), x(\lambda t), x(0)) + \int_0^b K_i(t, s, x(s), x(\lambda s)) ds, \quad t \in [0, b], \quad 0 < \lambda < 1.$$

The operator  $A_2$  is an increasing operator,  $A_1, A_2, A_3$  are Picard operators and  $A_1 \leq A_2 \leq A_3$ . By Lemma 2.2 it follows that  $x_1^* \leq x_2^* \leq x_3^*$ .



## 4 Numerical examples

**Example 4.1** Consider the following equation:

$$x(t) = t + \frac{1}{200}x(t) + \frac{1}{160}x\left(\frac{t}{2}\right) + 791 + \int_0^3 t(t+s + \frac{1}{45}\cos(x(s)) + \frac{1}{60}\sin(x(\frac{s}{2})))ds, \quad t \in [0, 3] \quad (4.1)$$

In this case  $X := \mathbb{R}$ ,  $b := 3$ ,  $\lambda := \frac{1}{2}$ ,  $g(t, u, v, \beta) := t + \frac{1}{200}u + \frac{1}{160}v + 791$ ,  $K(t, s, u, v) := t(t + s + \frac{1}{45}\cos u + \frac{1}{60}\sin v)$ ,  $L_g := \frac{1}{160}$ ,  $L_K := \frac{1}{15}$ ,  $b := 3$ ,  $\alpha^* := 800$ .

We have

**Theorem 4.2** The equation (4.1) has a unique solution in  $C[0, 3]$ .

**Example 4.3** Consider the equation

$$x(t) = 2t^2 + \frac{1}{5}\sin x(t) + \frac{1}{5}x\left(\frac{t}{4}\right) + \frac{4}{5}x(0) - \frac{1}{5} + \int_0^5 t(t+s + \frac{1}{150}\cos x(s) + \frac{1}{100}\sin(x(\frac{s}{4})))ds, \quad t \in [0, 5] \quad (4.2)$$

Here  $X := \mathbb{R}$ ,  $b := 5$ ,  $\lambda := \frac{1}{4}$ ,  $g(t, u, v, \beta) := 2t^2 + \frac{1}{5}\sin u + \frac{1}{5}v + \frac{4}{5}\beta - \frac{1}{5}$ ,  $K(t, s, u, v) := t(t + s + \frac{1}{150}\cos u + \frac{1}{100}\sin v)$ ,  $L_g := \frac{1}{5}$ ,  $L_K := \frac{1}{20}$ ,  $b := 5$ .

For (4.2) the equation  $\alpha = g(0, \alpha, \alpha, \alpha)$  becomes  $\sin \alpha = 1$ . By using Theorem 3.1 we obtain:

**Theorem 4.4** The equation (4.2) has solutions in  $C[0, 5]$ .

**Example 4.5** Consider the equation

$$x(t) = t + \frac{2}{3}\cos x(t) + \frac{1}{3}x\left(\frac{t}{3}\right) + \frac{2}{3}x(0) - 2 + \int_0^2 t(t+s + \frac{1}{8}\cos(x(s)) + \frac{1}{24}\sin(x(\frac{s}{3})))ds, \quad t \in [0, 2]$$

In this case  $X := \mathbb{R}$ ,  $b := 2$ ,  $\lambda := \frac{1}{3}$ ,  $g(t, u, v, \beta) := t + \frac{2}{3}\cos u + \frac{1}{3}v + \frac{2}{3}\beta - 2$ ,  $K(t, s, u, v) := t(t + s + \frac{1}{8}\cos u + \frac{1}{24}\sin v)$ ,  $L_g := \frac{2}{3}$ ,  $L_K := \frac{1}{4}$ .

By considering  $g(0, \alpha, \alpha, \alpha) = \alpha$  we obtain  $\cos \alpha = 3$  and  $S_g = \emptyset$ .

So, for this equation we can't apply the previous results.

**Remark 4.6** In the papers [16] and [17] we have considered some Volterra integral equations and we applied the weakly Picard operators technique to obtain existence and data dependence results for the solutions of those equations.

## References

- [1] N.V. AZBELEV, V.P. MAKSIMOV, L.F. RAHMATULINA, *Introduction to functional-differential equations theory*, MIR, Moscow, 1991 (in Russian)
- [2] R.E. BELLMAN, K.L. COOKE, *Differential difference equations*, Acad.Press, New York, 1963.
- [3] S.R. BERNFELD, V. LAKSHMIKANTHAM, *An introduction to nonlinear boundary value problems*, Acad.Press, New York, 1974.
- [4] C. CORDUNEANU, *Integral equations and applications*, Cambridge University Press, 1991.
- [5] L.F. ELSGOLTZ, S.B. NORKIN, *Introduction to the theory of differential equations with deviating arguments*, MIR, Moscow, 1971 (in Russian)
- [6] K. GOPALSAMY, *Stability and oscillations in delay differential equations of population dynamics*, Kluwer, Academic Publisher, 1992.
- [7] D. GUO, V. LAKSHMIKANTHAM, X. LIU, *Nonlinear integral equations in abstracts spaces*, Kluwer, Dordrecht, 1996.
- [8] J.K. HALE, *Theory of functional-differential equations*, Springer Verlag, 1977.
- [9] J.K. HALE, SJOERD M. VERDUYN LUNEL *Introduction to functional-differential equations*, Springer Verlag, New York, 1993.
- [10] A. HALANAY, *Differential equations: stability, oscillations, time lags*, Academic Press, New York, 1965.
- [11] A. ISERLES, On the generalized pantograph functional-differential equation, *European J.Appl.Math.* (1992), 1-38.
- [12] T. KATO, J.B. MCLEOD, The functional-differential equation  $y'(x) = ay(\lambda x) + by(x)$ , *Bull.Amer.Math.Soc.*, 77, 6(1971), 891-937.
- [13] V. LAKSHMIKANTHAM (ed), *Trends in the theory and practice of nonlinear differential equations*, Marcel Dekker, New York, 1984.
- [14] V. MURESAN, *Differential equations with affine modification of the argument*, Transilvania Press, Cluj-Napoca, 1997 (in Romanian)
- [15] V. MURESAN, *Functional-integral equations*, Ed. Mediamira, Cluj-Napoca, 2003.
- [16] V. MURESAN, A functional integral equation with linear modification of the argument, via weakly Picard operators, *International Conference on Nonlinear Operators, Differential Equations and Applications (ICNODEA)*, July 4-8, 2007, Cluj-Napoca, Romania, Fixed Point Theory, vol. 9(2008), no.1, 189-197.
- [17] V. MURESAN, Weakly Picard operators' technique, applied to a class of functional-integral equations, *Proceedings of the 7-th International Conference APLIMAT 2008*, February 5-8, 2008, Bratislava, 257-264, APLIMAT - Journal of Applied Mathematics, vol. 1(2008), no.1, 185-194.
- [18] J.R. OCKENDON, Differential equations and industry, *The Math. Scientist.* 5(1980) no.1, 1-12.
- [19] W. POGORZELSKI, *Integral equations and their applications*, Pergamon Press, 1966.
- [20] R. PRECUP, *Integral and nonlinear equations*, Babes-Bolyai University, Cluj-Napoca, 1993 (in Romanian).

- [21] I.A. RUS, *Metrical fixed point theorems*, Univ. of Cluj-Napoca, 1979.
- [22] I.A. RUS, *Weakly Picard mappings*, Comment. Math. Univ. Caroline, 34, 4(1993), 769-773.
- [23] I.A. RUS, Picard operators and applications, *Scientia Math. Japonicae*, 58(2003), Nr.1, 191-219.
- [24] I.A. RUS, *Weakly Picard operators and applications*, Seminar on Fixed Point Theory, 2(2001), 41-58.
- [25] I. A. Rus, *Generalized contractions and applications*, Cluj Univ. Press., Cluj-Napoca, 2001.
- [26] I.A. RUS, A class of nonlinear functional-integral equations, via weakly Picard operators, *Anal. Univ. Craiova, ser. Mat-Inf.*, 28(2001), 10-15.
- [27] A. SINCELEAN, On a class of functional-integral equations, *Seminar on Fixed Point Theory*, 1 (2000), 87-92.

### Current address

#### Muresan Viorica, Professor

Department of Mathematics, Faculty of Automation and Computer Sciences,  
Technical University of Cluj-Napoca, 25, C. Daicoviciu Street, Cluj-Napoca, Romania  
e-mail address: vmuresan@math.utcluj.ro



ON ASYMPTOTIC BEHAVIOUR  
OF THE NONLINEAR VISCOELASTIC  
MINDLIN-TIMOSHENKO THIN PLATE MODEL

PANCZA Dávid, (SK)

**Abstract.** The oscillation of a thin viscoelastic plate is described by a time dependent vector function. Under a constant load of forces the oscillation has to relax to a final constant state. If the mathematical model of the thin plate is correct, the convergence of the force function to a constant must imply a convergence of the state function to a stationary vector. In our paper we are testing the relaxing of a nonlinear Mindlin-Timoshenko thin plate model.

**Key words and phrases.** Mindlin-Timoshenko thin plate model, viscoelasticity, material function, asymptotic behaviour.

*Mathematics Subject Classification.* Primary 74D10; Secondary 35Q72.

## 1 Introduction

The formulation of the viscoelastic Mindlin-Timoshenko (MT) thin plate model consists of three differential equations (Lagnese – Lions, 1989) with additional convolution terms. In the case of larger deformations some nonlinearities arise. We shall deal with a simplified model with nonlinearities only in the third coordinate. Our aim is to test the behaviour of its weak solution for  $t \rightarrow \infty$ . Convergence of the solution to a stationary vector in the case of a constant load was already proved for a Kirchhoff model with an exponentially decreasing material function (I. Bock, 2002). We apply a similar method to our model admitting more general material functions.

## 2 The formulation of the problem

Consider a region  $\Omega \in \mathbb{R}$  with a Lipschitz boundary  $\Gamma = \Gamma_0 \cup \Gamma_1$ ,  $\Gamma_0 \neq \emptyset$ . We introduce the following Sobolev spaces and norms:

$$W^{1,p}(\Omega) = \{v \in L^p(\Omega); \partial_\alpha v \in L^p(\Omega), \alpha = 1, 2, p > 1\}$$

$$\|v\|_{1,p} = \left[ \iint_{\Omega} v^p + (\partial_1 v)^p + (\partial_2 v)^p d\mathbf{x} \right]^{1/p},$$

and

$$W_{\Gamma_0}^{1,p}(\Omega) = \{v \in W^{1,p}(\Omega); v = 0 \text{ on } \Gamma_0\},$$

$$\|v\|_{1,p,0} = \left[ \iint_{\Omega} (\partial_1 v)^p + (\partial_2 v)^p d\mathbf{x} \right]^{1/p}.$$

The partial derivatives are considered in the sense of distributions and the boundary condition in the sense of traces in the space  $W^{1,p}(\Omega)$ . Moreover we denote

$$\mathcal{W} = (W_{\Gamma_0}^{1,2}(\Omega))^3, \quad \|\mathbf{y}\|_{\mathcal{W}} = \sum_{i=1}^3 \|y_i\|_{1,2,0},$$

$$\mathcal{V} = (W_{\Gamma_0}^{1,2}(\Omega))^2 \times W_{\Gamma_0}^{1,4}(\Omega), \quad \|\mathbf{y}\|_{\mathcal{V}} = \|y_1\|_{1,2,0} + \|y_2\|_{1,2,0} + \|y_3\|_{1,4,0}.$$

Let  $\mathbf{w} = (\phi_1, \phi_2, w)$  and  $\mathbf{v} = (\psi_1, \phi_2, v)$  be vectors from the space  $\mathcal{W}$ . Consider an operator  $\mathcal{A}_2 : \mathcal{W} \rightarrow \mathcal{W}^*$  satisfying the following conditions:

a)  $\langle \mathcal{A}_2(\mathbf{w}), \mathbf{v} \rangle$  is a bilinear form that can be written as a sum

$$\langle \mathcal{A}_2(\mathbf{w}), \mathbf{v} \rangle = \iint_{\Omega} \sum_{i=1}^n L_i(\mathbf{w}) L_i(\mathbf{v}) dx$$

where  $L_i : \mathcal{W} \rightarrow L^2(\Omega)$  are linear operators,

b) there exists a constant  $c_2 > 0$  such that

$$\langle \mathcal{A}_2(\mathbf{w}), \mathbf{w} \rangle \geq c_2 (\|\phi_1\|_{1,2,0}^2 + \|\phi_2\|_{1,2,0}^2 + \|w\|_{1,2,0}^2) \geq \frac{c_2}{3} \|\mathbf{w}\|_{\mathcal{W}}^2,$$

c) there exists a constant  $C_2 > 0$  such that

$$\langle \mathcal{A}_2(\mathbf{w}), \mathbf{v} \rangle \leq C_2 \|\mathbf{w}\|_{\mathcal{W}} \|\mathbf{v}\|_{\mathcal{W}}.$$

For  $\mathbf{w} = (\phi_1, \phi_2, w)$ ,  $\mathbf{v} = (\psi_1, \psi_2, w)$  from the space  $\mathcal{V}$  we define an operator  $\mathcal{A}_4 : \mathcal{V} \rightarrow \mathcal{V}^*$ :

$$\langle \mathcal{A}_4(\mathbf{w}), \mathbf{v} \rangle = \langle \mathcal{A}_4(w), v \rangle = \frac{1}{2} \iint_{\Omega} (\nabla w \cdot \nabla w)(\nabla w \cdot \nabla v) d\mathbf{x}.$$

Finally, let there be given a bounded linear operator  $\mathcal{F} \in C^1([0, \infty), \mathcal{W}^*)$ :

$$\langle \mathcal{F}(t), \mathbf{v} \rangle \leq K_f \|\mathbf{v}\|_{\mathcal{W}}, \quad K_f > 0.$$

Let  $D \in C^1([0, \infty), \mathbb{R}^+)$  be a function with  $\int_0^\infty |D'(t)| dt \leq M_d$  where  $M_d > 0$  is a constant. Let  $a > D(0)$  be a constant too.

The weak formulation of our problem is an operator equation (1) in  $\mathcal{V}^*$  with one boundary condition (2):

$$a\mathcal{A}_2(\mathbf{w}) + D' * \mathcal{A}_2(\mathbf{w}) + \mathcal{A}_4(\mathbf{w}) = \mathcal{F} \quad \forall t \in [0, \infty), \quad (1)$$

$$\mathbf{w} = \mathbf{0} \text{ on } \Gamma_0. \quad (2)$$

The system (1-2) has a unique bounded solution  $\mathbf{w} \in C([0, T], \mathcal{V})$  (Pancza, 2002).

Our aim is to investigate the situation, when  $\mathcal{F}$  tends to a constant in time. We shall prove that the solution  $\mathbf{w}$  of (1-2) asymptotically converges to a constant in time function.

### 3 Properties of the operators $\mathcal{A}_2, \mathcal{A}_4$

From now on we suppose that for all  $t \in [0, \infty)$  the material function  $D(t) \in C^2[0, \infty)$  satisfies the following conditions with constants  $\beta > 0$ ,  $\alpha_i > 0$  for  $i = 1, 3, 5$  and  $\alpha_i \geq 0$  for  $i = 2, 4, 6$ :

$$\alpha_2 e^{-\beta t} \leq D(t) \leq \alpha_1 e^{-\beta t}, \quad -\alpha_3 e^{-\beta t} \leq D'(t) \leq -\alpha_4 e^{-\beta t}, \quad \alpha_6 e^{-\beta t} \leq D''(t) \leq \alpha_5 e^{-\beta t}.$$

We denote  $d_0 = D(0)$ ,  $d_1 = D'(0)$  and  $d_2 = D''(0)$ . For  $t = 0$  we get the basic inequalities that must be satisfied by the constants:

$$\alpha_6 \leq d_2 \leq \alpha_5, \quad \frac{\alpha_6}{\beta} \leq \alpha_4 \leq -d_1 \leq \alpha_3 \leq \frac{\alpha_5}{\beta}, \quad \frac{\alpha_4}{\beta} \leq \alpha_2 \leq d_0 \leq \alpha_1 \leq \frac{\alpha_3}{\beta}$$

**Lemma 3.1** *Let  $\beta > 0$ ,  $0 \leq \gamma < 2\beta$  be constants. We denote*

$$P(t) = \int_0^t \left\langle \int_0^s e^{(\gamma-\beta)s+\beta\tau} \mathcal{A}_2(\mathbf{w}(\tau)) d\tau, \mathbf{w}(s) \right\rangle ds, \quad (3)$$

$$l_i(\mathbf{w}) = \|L_i(\mathbf{w})\|_{L^2(\Omega)}, \quad M_i = \int_0^t e^{\gamma s} l_i^2(\mathbf{w})(s) ds, \quad M = \sum_{i=1}^n M_i.$$

For all  $t \in [0, \infty)$  holds the estimate

$$|P| \leq \frac{2}{2\beta - \gamma} M. \quad (4)$$

**Proof.** We denote  $N_i(t) = \int_0^t \int_0^s e^{(\gamma-\beta)s+\beta\tau} l_i(\mathbf{w})(\tau) d\tau l_i(\mathbf{w})(s) ds$ . After applying the Cauchy – Schwarz inequality on (3) we obtain

$$|P(t)| \leq \sum_{i=1}^n N_i(t). \quad (5)$$

Using the same inequality once again we have

$$N_i \leq \Lambda_i^{\frac{1}{2}} M_i^{\frac{1}{2}}, \quad (6)$$

where

$$\Lambda_i = \int_0^t \left( e^{(\frac{\gamma}{2}-\beta)s} \int_0^s e^{\beta\tau} l_i(\mathbf{w})(\tau) d\tau \right)^2 ds.$$

After integrating by parts the terms  $\Lambda_i$  we have:

$$\Lambda_i = -\frac{e^{\gamma t}}{2\beta - \gamma} \left( \int_0^t e^{\beta(t-\tau)} l_i(\mathbf{w})(\tau) d\tau \right)^2 + \frac{2}{2\beta - \gamma} N_i.$$

Hence we get the estimate

$$\Lambda_i \leq \frac{2}{2\beta - \gamma} N_i. \quad (7)$$

From (6), (7) we have  $N_i \leq \frac{2}{2\beta - \gamma} M_i$ . Applied into (5) we get (4).  $\square$

For  $\mathcal{A}_4$  it holds:

$$\begin{aligned} & \langle \mathcal{A}_4(\mathbf{w}) - \mathcal{A}_4(\mathbf{y}), \mathbf{w} - \mathbf{y} \rangle = \\ & \frac{1}{2} \iint_{\Omega} (\nabla w \cdot \nabla w)(\nabla w \cdot (\nabla w - \nabla y)) + (\nabla y \cdot \nabla y)(\nabla y \cdot (\nabla y - \nabla w)) d\mathbf{x} = \\ & = \frac{1}{4} \iint_{\Omega} (\nabla w \cdot \nabla w - \nabla y \cdot \nabla y)^2 + (\nabla w \cdot \nabla w + \nabla y \cdot \nabla y)((\nabla w - \nabla y) \cdot (\nabla w - \nabla y)) d\mathbf{x}. \end{aligned}$$

We define an operator  $\mathcal{B}_4$ :

$$\mathcal{B}_4(w, y) = \frac{1}{4} \iint_{\Omega} \frac{1}{2} (\nabla w \cdot \nabla w - \nabla y \cdot \nabla y)^2 + (\nabla w \cdot \nabla w)(\nabla(w - y) \cdot \nabla(w - y)) d\mathbf{x}. \quad (8)$$

**Lemma 3.2** *Let  $\mathbf{y} \in \mathcal{V}$  be a function constant in time and  $y$  its third coordinate. Then*

$$\langle \mathcal{A}_4(\mathbf{w}) - \mathcal{A}_4(\mathbf{y}), \mathbf{w} - \mathbf{y} \rangle \geq \mathcal{B}_4(w, y) \geq 0. \quad (9)$$

Moreover, if  $\mathbf{w}$  is a solution of (1-2), it holds:

$$\frac{d}{dt} \left( \frac{a}{2} \langle \mathcal{A}_2(\mathbf{u}), \mathbf{u} \rangle + \mathcal{B}_4(w, y) \right) (t) = \left\langle \frac{d}{dt} \mathfrak{A}(\mathbf{w}), \mathbf{w} - \mathbf{y} \right\rangle (t), \quad (10)$$

where  $\mathfrak{A} = a\mathcal{A}_2 + \mathcal{A}_4$ .

**Proof.** The inequality (9) is obvious. Before the proof of (10) we need an auxiliary estimate. We denote  $w = w(t)$ ,  $W = w(t + \eta)$ ,  $\mathbf{w} = \mathbf{w}(t)$ ,  $\mathbf{W} = \mathbf{w}(t + \eta)$  and

$$\mathcal{Z}(W, w) = \frac{1}{4} \iint_{\Omega} [(\nabla W - \nabla w) \cdot (\nabla W - \nabla w)][\nabla W \cdot \nabla W - \nabla w \cdot \nabla w] d\mathbf{x}.$$



It is not complicated to prove that

$$\mathcal{Z}(W, w) \leq \langle \mathcal{A}_4(W) - \mathcal{A}_4(w), W - w \rangle. \quad (11)$$

We subtract the equation (1) in time  $t$  from the same equation in time  $t + \eta$ . We apply the result on  $\mathbf{v} = \mathbf{W} - \mathbf{w}$ . Using the properties of  $\mathcal{A}_2$  we get

$$\begin{aligned} \langle \mathcal{A}_4(W) - \mathcal{A}_4(w), W - w \rangle &\leq \\ &\leq \langle \mathcal{F}(t + \eta) - \mathcal{F}(t) - D' * \mathcal{A}_2(\mathbf{W} - \mathbf{w}), \mathbf{W} - \mathbf{w} \rangle. \end{aligned} \quad (12)$$

Applying (12) into (11) and dividing the inequality by  $\eta$  we get for  $\eta \rightarrow 0$ :

$$\lim_{\eta \rightarrow 0} \frac{1}{\eta} \mathcal{Z}(W, w) \leq \left\langle \mathcal{F}' - d_1 \mathcal{A}_2(\mathbf{w}) - D'' * \mathcal{A}_2(\mathbf{w}), \lim_{\eta \rightarrow 0} (\mathbf{W} - \mathbf{w}) \right\rangle = 0.$$

Now we can prove (10). Let  $\mathbf{u}(t) = \mathbf{w}(t) - \mathbf{y}$ ,  $\mathbf{U}(t) = \mathbf{u}(t + \eta)$  and let  $u, U$  be their third coordinates. We have

$$\begin{aligned} \langle \mathcal{A}_4(W) - \mathcal{A}_4(w), \frac{1}{2}(U + u) \rangle &= \\ &= \frac{1}{4} \iint_{\Omega} (\nabla W \cdot \nabla W)(\nabla W \cdot \nabla U + \nabla u) - (\nabla w \cdot \nabla w)(\nabla w \cdot \nabla U + \nabla u) d\mathbf{x} = \\ &= \frac{1}{8} \iint_{\Omega} Q(\nabla W + \nabla w) \cdot (\nabla U + \nabla u) + R(\nabla W - \nabla w) \cdot (\nabla U + \nabla u) d\mathbf{x}, \end{aligned}$$

where

$$Q = \nabla W \cdot \nabla W - \nabla w \cdot \nabla w, \quad R = \nabla W \cdot \nabla W + \nabla w \cdot \nabla w.$$

Using the same notation we can write:

$$\begin{aligned} \mathcal{B}_4(W, y) - \mathcal{B}_4(w, y) &= \\ &= \frac{1}{8} \iint_{\Omega} 2Q(\nabla W \cdot \nabla U + \nabla w \cdot \nabla u) + R((\nabla W - \nabla w) \cdot (\nabla U + \nabla u)) d\mathbf{x}. \end{aligned}$$

Comparing the results we see that

$$\mathcal{B}_4(W, y) - \mathcal{B}_4(w, y) = \left\langle \mathcal{A}_4(W) - \mathcal{A}_4(w), \frac{1}{2}(U + u) \right\rangle + \mathcal{Z}(W, w). \quad (13)$$

Using (13) and the bilinearity of  $\mathcal{A}_2$  we get the equation

$$\begin{aligned} \frac{1}{\eta} \left( \langle a \mathcal{A}_2(\mathbf{U}) + \mathcal{A}_4(W) - a \mathcal{A}_2(\mathbf{u}) - \mathcal{A}_4(w), \frac{1}{2}(U + u) \rangle + \mathcal{Z}(W, w) \right) &= \\ &= \frac{1}{\eta} \left( \frac{a}{2} (\langle \mathcal{A}_2(\mathbf{U}), \mathbf{U} \rangle - \langle \mathcal{A}_2(\mathbf{u}), \mathbf{u} \rangle) + (\mathcal{B}_4(W, y) - \mathcal{B}_4(w, y)) \right). \end{aligned} \quad (14)$$

For  $\eta \rightarrow 0$  the term  $\mathcal{Z}$  tends to zero and the rest of the terms on the left-hand side define the first derivative of the operator  $\mathfrak{A}(\mathbf{w})$  applied on  $\mathbf{u}$ . This derivative can be evaluated if  $\mathbf{w}$  is a solution of (1). So the limit of the left-hand side of (14) exists, hence the limit of the right-hand side must exist too and it holds:

$$\left\langle \frac{d}{dt}(\mathfrak{A}(\mathbf{w}), \mathbf{u}) = \frac{d}{dt} \left( \frac{a}{2} \langle \mathcal{A}_2(\mathbf{u}), \mathbf{u} \rangle + \mathcal{B}_4(w, y) \right) \right.$$

□

**Lemma 3.3** *Let  $\chi \in C([0, \infty))$  and  $\chi(t) \rightarrow 0$  for  $t \rightarrow \infty$ . Then for  $\gamma > 0$  it holds that*

$$\int_0^t e^{-\gamma(t-s)} \chi(s) ds \rightarrow 0 .$$

**Proof.** Let there be given  $\epsilon > 0$ . We set  $\eta = \frac{\epsilon\gamma}{2}$ . From the suppositions follows:

- a)  $\forall \eta > 0 : \exists t_\eta, \forall t > t_\eta : \chi(t) < \eta,$
- b)  $\exists M > 0, \forall t \in [0, \infty) : |\chi(t)| \leq M.$

The integral can be divided in two parts and bounded in the following way:

$$\begin{aligned} \int_0^t e^{-\gamma(t-s)} \chi(s) ds &\leq M \int_0^{t_\eta} e^{-\gamma(t-s)} ds + \eta \int_{t_\eta}^t e^{-\gamma(t-s)} ds \leq \\ &\leq \frac{M}{\gamma} (e^{-\gamma(t-t_\eta)} - e^{-\gamma t}) + \eta \gamma \leq \frac{M}{\gamma} e^{-\gamma(t-t_\eta)} + \frac{\epsilon}{2} . \end{aligned}$$

The last term is for every  $t > t_\epsilon$ , where

$$t_\epsilon = t_\eta + \frac{1}{\gamma} \ln \left( \frac{2M}{\epsilon\gamma} \right) ,$$

less than  $\epsilon$ , and that proves the lemma.

□

#### 4 The asymptotic behaviour of the model for $t \rightarrow \infty$ at constant load

We suppose that  $\mathcal{F} \in C^1([0, \infty), \mathcal{W}^*)$  tends for  $t \rightarrow \infty$  to a constant load  $\mathcal{F}_\infty$ , i.e.

$$\lim_{t \rightarrow \infty} \|\mathcal{F}(t) - \mathcal{F}_\infty\|_{\mathcal{W}^*} = 0, \quad \lim_{t \rightarrow \infty} \|\mathcal{F}'(t)\|_{\mathcal{W}^*} = 0. \quad (15)$$

Our hypothesis is that the material function  $D$  will relax to zero and the system will converge to a stationary state, i.e.  $\mathbf{w}(t) \rightarrow \mathbf{w}_\infty$  for  $t \rightarrow \infty$ , where  $\mathbf{w}_\infty$  represents the final state. The convolution term of the equation must converge to a constant vector

$$\int_0^t D'(t-\tau) \mathcal{A}_2(\mathbf{w}(\tau)) d\tau \rightarrow -d_0 \mathcal{A}_2(\mathbf{w}_\infty)$$

and the original equation (1) to a stationary equation

$$(a - d_0)\mathcal{A}_2(\mathbf{w}) + \mathcal{A}_4(w) = \mathcal{F}_\infty. \quad (16)$$

The final state  $\mathbf{w}_\infty$  has to be a solution of the system (16). The problem (16) is in fact the problem (1) without convolution terms and the existence and uniqueness of its solution can be proven by the Lax-Milgram theorem (Lagnese-Lions, 1989).

Let us denote  $\mathbf{w}_\infty$  the solution of (16) and  $\mathbf{u} = \mathbf{w} - \mathbf{w}_\infty$ . In order to get an estimate of  $\mathbf{u}$  we need to reformulate the original system (1). We carry out the first derivative of (1):

$$D'' * \mathcal{A}_2(\mathbf{w}) + d_1 \mathcal{A}_2(\mathbf{w}) + \frac{d}{dt}[a \mathcal{A}_2(\mathbf{w}) + \mathcal{A}_4(\mathbf{w})] = \mathcal{F}'$$

and add to it the same equation (1) multiplied by  $\beta$ . We get

$$(\beta D' + D'') * \mathcal{A}_2(\mathbf{w}) + (d_1 + \beta a) \mathcal{A}_2(\mathbf{w}) + \beta \mathcal{A}_4(\mathbf{w}) + \frac{d}{dt}[a \mathcal{A}_2(\mathbf{w}) + \mathcal{A}_4(\mathbf{w})] = \beta \mathcal{F} + \mathcal{F}' \quad (17)$$

$$\mathbf{w}(0) = \mathfrak{A}^{-1}(\mathcal{F}(0)). \quad (18)$$

Now we can subtract the equation (16) multiplied by  $\beta$  from (17). We use that it holds:

$$\mathcal{A}_2(\mathbf{w}) = \mathcal{A}_2(\mathbf{u}) + \mathcal{A}_2(\mathbf{w}_\infty)$$

$$G' * \mathcal{A}_2(\mathbf{w})(t) = G' * \mathcal{A}_2(\mathbf{u})(t) + (G(t) - G(0)) \mathcal{A}_2(\mathbf{w}_\infty),$$

where  $G$  stands for  $\beta D + D'$ . The result is

$$(\beta D' + D'') * \mathcal{A}_2(\mathbf{u}) + (\beta a + d_1) \mathcal{A}_2(\mathbf{u}) + \beta (\mathcal{A}_4(\mathbf{w}) - \mathcal{A}_4(\mathbf{w}_\infty)) + \frac{d}{dt}[a \mathcal{A}_2(\mathbf{u}) + \mathcal{A}_4(\mathbf{w})] = \mathcal{H}, \quad (19)$$

$$\text{where } \mathcal{H} = \beta (\mathcal{F} - \mathcal{F}_\infty) + \mathcal{F}' - (\beta D + D') \mathcal{A}_2(\mathbf{w}_\infty).$$

According to suppositions (15) it holds that  $\mathcal{H} \rightarrow 0$  for  $t \rightarrow \infty$  in the norm of the space  $\mathcal{W}^*$ .

From (19) we are going now to extract an estimate of  $\mathbf{u}$  in order to prove that  $\mathbf{u} \rightarrow 0$  in the norm of an appropriate space. Let us apply the operators of (19) on the vector  $\mathbf{u}$  and multiply the whole equation by  $\exp(\gamma t)$  (the constant  $\gamma$ ,  $0 < \gamma < 2\beta$  will be determined later). After carrying out the integration in the time variable  $t$  we get an equation

$$I_1 + I_2 + I_3 = \int_0^t \langle \mathcal{H}, \mathbf{u} \rangle e^{\gamma s} ds, \quad (20)$$

$$\text{where } I_1 = \int_0^t (\beta a + d_1) \langle \mathcal{A}_2(\mathbf{u}), \mathbf{u} \rangle + \beta \langle \mathcal{A}_4(w) - \mathcal{A}_4(w_\infty), u \rangle e^{\gamma s} ds,$$

$$I_2 = \int_0^t \left( \left\langle \frac{d}{dt}[a \mathcal{A}_2(\mathbf{w}) + \mathcal{A}_4(w)], \mathbf{u} \right\rangle \right) e^{\gamma s} ds,$$

$$I_3 = \int_0^t \langle (\beta D' + D'') * \mathcal{A}_2(\mathbf{u}), \mathbf{u} \rangle e^{\gamma s} ds.$$

Using (10)  $I_2$  can be written in the form

$$I_2 = \int_0^t \left( \frac{d}{dt} \left[ \frac{a}{2} \langle \mathcal{A}_2(\mathbf{u}), u \rangle + \mathcal{B}_4(w, y) \right] \right) e^{\gamma s} ds.$$

We start with the estimate of  $I_3$ . From the properties of  $D$  it follows that

$$(-\beta\alpha_3 + \alpha_6)e^{-\beta t} \leq (\beta D' + D'')(t) \leq (-\beta\alpha_4 + \alpha_5)e^{-\beta t}. \quad (21)$$

The exponential function on the left-hand side is negative due to (3.2) and that on the right-hand side is positive. So we can express inequalities (21) in the following way:

$$|(\beta D' + D'')(t)| < \alpha e^{-\beta t}, \quad \text{where } \alpha = \max\{\beta\alpha_3 - \alpha_6, \alpha_5 - \beta\alpha_4\}.$$

Now using lemma 1 we get

$$|I_3| < \frac{2\alpha}{2\beta - \gamma} \int_0^t e^{\gamma s} \sum_{i=1}^n l_i^2(\mathbf{u}) ds. \quad (22)$$

The estimate of the sum  $I_1 + I_3$  is possible only at some additional conditions:

**Lemma 4.1** *Let the constant  $a$  satisfy*

$$\frac{\alpha}{\beta^2} - \frac{d_1}{\beta} < a. \quad (23)$$

*Then there exist constants  $\theta > 0$ ,  $\gamma_0 > 0$  such that  $\forall \gamma: 0 < \gamma \leq \gamma_0$  it holds:*

$$I_1 + I_3 \geq J,$$

$$\text{where } J = \int_0^t (\theta \langle \mathcal{A}_2(\mathbf{u}), \mathbf{u} \rangle + \beta \langle \mathcal{A}_4(w) - \mathcal{A}_4(w_\infty), u \rangle) e^{\gamma s} ds.$$

**Proof.** Supposing (23) it is possible to find numbers  $\theta > 0$ ,  $\gamma_0 > 0$  such that  $\forall \gamma: 0 < \gamma \leq \gamma_0$  it holds

$$\theta + \frac{2\alpha}{2\beta - \gamma} \leq \theta + \frac{2\alpha}{2\beta - \gamma_0} < \beta a + d_1.$$

Hence it follows

$$I_1 \geq \int_0^t \left( \left( \theta + \frac{2\alpha}{2\beta - \gamma} \right) \langle \mathcal{A}_2(\mathbf{u}), \mathbf{u} \rangle + \beta \langle \mathcal{A}_4(\mathbf{w}) - \mathcal{A}_4(\mathbf{w}_\infty), \mathbf{u} \rangle \right) e^{\gamma s} ds.$$

Using (4) and (22) we obtain

$$I_1 + I_3 \geq I_1 - |I_3| \geq J$$

Let us continue with the estimate of the term  $J$ . We set

$$\gamma = \min\left\{\beta, \gamma_0, \frac{2\theta}{a}\right\}.$$

Using (9) we can bound the term containing  $\mathcal{A}_4$ :

$$J \geq \int_0^t \gamma \left( \frac{a}{2} \langle \mathcal{A}_2(\mathbf{u}), \mathbf{u} \rangle + \mathcal{B}_4(w, w_\infty) \right) e^{\gamma s} ds =: J_1.$$

Now we are able to carry out the integration in  $t$  of  $J_1 + I_2$ :

$$J_1 + I_2 = \left( \frac{a}{2} \langle \mathcal{A}_2(\mathbf{u}), \mathbf{u} \rangle(t) + \mathcal{B}_4(\mathbf{w}, \mathbf{w}_\infty)(t) \right) e^{\gamma t} - K_2 - K_4, \quad (24)$$

where  $K_2 = \frac{a}{2} \langle \mathcal{A}_2(\mathbf{u}(0)), \mathbf{u}(0) \rangle$  and  $K_4 = \mathcal{B}_4(\mathbf{w}(0), \mathbf{w}_\infty)$ .

According to properties of  $\mathcal{A}_2$  we can write

$$J + I_2 \geq \frac{ac_2}{6} \|\mathbf{u}\|_{\mathcal{W}}^2 e^{\gamma t} - K_2 - K_4. \quad (25)$$

Coming back to (20) and introducing there the obtained estimates we get

$$\frac{ac_2}{6} \|\mathbf{u}\|_{\mathcal{W}}^2 \leq \int_0^t \|\mathcal{H}\|_{\mathcal{W}^*} \|\mathbf{u}\|_{\mathcal{W}} e^{-\gamma(t-s)} ds + (K_2 + K_4) e^{-\gamma t}. \quad (26)$$

The solutions of (1-2) are bounded, so we can a priori bound the unknown vector  $\mathbf{u} = \mathbf{w} - \mathbf{w}_\infty$  with an appropriate constant  $K_u$ . Hence according to lemma 3 it holds for  $t \rightarrow \infty$ :

$$K_u \int_0^t \|\mathcal{H}\|_{\mathcal{W}^*} e^{\gamma(t-s)} ds \rightarrow 0,$$

so  $\|\mathbf{u}\|_{\mathcal{W}} \rightarrow 0$ . This implies  $\mathbf{w} \rightarrow \mathbf{w}_\infty$  in the space  $\mathcal{W}$ .

Using this result we can prove the convergence of the last coordinate  $u = w - w_\infty$  in the norm of the space  $W_{\Gamma_0}^{1,4}$ . Let us subtract (16) from (1) and apply the result on  $u$ :

$$\begin{aligned} (a - d_0) \langle \mathcal{A}_2(\mathbf{u}), \mathbf{u} \rangle + \langle \mathcal{A}_4(\mathbf{w}) - \mathcal{A}_4(\mathbf{w}_\infty), \mathbf{u} \rangle &= \\ = \langle \mathcal{F} - \mathcal{F}_\infty, \mathbf{u} \rangle - \langle d_0 \mathcal{A}_2(\mathbf{w}) - D' * \mathcal{A}_2(\mathbf{w}), \mathbf{u} \rangle. \end{aligned} \quad (27)$$

The solution  $\mathbf{w}$  is bounded, hence we get

$$\|\langle d_0 \mathcal{A}_2(\mathbf{w}) \rangle_{\mathcal{W}^*}\| \leq d_0 C_2 M_w,$$

$$\|D' * \mathcal{A}_2(\mathbf{w})\|_{\mathcal{W}^*} \leq \alpha_1 C_2 M_w.$$

The left-hand side of (27) can be reduced using lemma 3 and we obtain the inequality

$$\|u\|_{1,4,0}^4 \leq (\|\langle \mathcal{F} - \mathcal{F}_\infty \rangle_{\mathcal{W}^*}\| + (d_0 + \alpha_1) C_2 M_w) \|\mathbf{u}\|_{\mathcal{W}} \rightarrow 0.$$

This implies the convergence of  $u$ .

**Theorem 4.2** *Let  $D$  and  $a$  satisfy suppositions (1-2), (23). Let  $\mathcal{F}$  satisfy (15). Then for  $t \rightarrow \infty$  the solution  $\mathbf{w}$  of (1-2) converges in the norm of  $\mathcal{V}$  to a stationary solution  $\mathbf{w}_\infty$  of (16).*

## Acknowledgement

The author gratefully acknowledges the Scientific Grant Agency VEGA for supporting this work under the Grant No. 1/4214/07.

**References**

- [1] I. BOCK: On the generalized von Kármán system system for viscoelastic plates, II. Long memory model, Department of Mathematics preprint series, FEI STU, Bratislava, 2002
- [2] J.E. LAGNESE, J.-L. LIONS: Modelling analysis and control of thin plates, Masson, Springer-Verlag, Paris, Berlin, 1989
- [3] D. PANCZA: The Solution of a Simplified Nonlinear Viscoelastic Mindlin-Timoshenko Thin Plate Model, Journal of Electrical Engineering, 53, 2002

**Current address**

**David Pancza**

Department of Mathematics, Faculty of Electrical Engineering and Information Technology,  
Slovak University of Technology,  
Ilkovičova 3, 812 19 Bratislava 1, Slovak Republic; tel.number +421260291617,  
e-mail: david.panczastuba.sk

## GENERALIZATION OF CERTAIN INTEGRAL INEQUALITIES

ŠMARDA Zdeněk, (CZ)

**Abstract.** In the paper certain integral inequalities are generalized . There are established conditions of more precise bounds of these inequalities and results are applied to investigation of boundedness of solutions of nonlinear integrodifferential equations .

**Key words and phrases.** Integral inequalities, boundedness of solutions, integrodifferential equations.

*Mathematics Subject Classification.* 45J05.

## 1 Introduction

Integral inequalities play a significant role in the study of differential, integral and integrodifferential equations (see[1-9]). For example, in the theory of differential and integrodifferential equations the integral inequalities are used to the study of boundedness and stability of solutions (see [8,9]). In [5] J.A. Oguntuase tried to obtain the generalizations of the Pachpatte's inequalities [6] but his proofs and assumptions were incorrect. The aim of the paper is to correct his results and also obtain a bound of the general version of inequalities in [5].

## 2 Main results

**Theorem 2.1** *Suppose that the functions  $u(t), f(t) \in C[I, R_+]$ ,  $k(t, s) \in C[I \times I, R_+]$ ,  $k_t(t, s) \in C[I \times I, R_-]$ ,  $I = [a, b]$  ,  $c$  be a nonnegative constant. If the inequality*

$$u(t) \leq c + \int_a^t f(s)u(s)ds + \int_a^t f(s) \left( \int_a^s k(s, \tau)u(\tau)d\tau \right) ds, \quad a \leq \tau \leq s \leq t \leq b \quad (1)$$

holds then

$$u(t) \leq c \left[ 1 + \int_a^t f(s) \exp \left( \int_a^s (f(\tau) + k(\tau, \tau)) d\tau \right) ds \right]. \quad (2)$$

**Proof.** Define a function  $v(t)$  by the right hand side of (1) . Then it follows that

$$u(t) \leq v(t). \quad (3)$$

Therefore

$$v'(t) = f(t)u(t) + f(t) \int_a^t k(t, \tau)u(\tau)d\tau \leq f(t) \left( v(t) + \int_a^t k(t, \tau)v(\tau)d\tau \right). \quad (4)$$

If we put

$$m(t) = v(t) + \int_a^t k(t, \tau)v(\tau)d\tau, \quad (5)$$

then it is obvious that  $v(t) \leq m(t)$ ,  $m(a) = v(a) = c$ . Thus

$$\begin{aligned} m'(t) &= v'(t) + k(t, t)v(t) + \int_a^t k_t(t, \tau)v(\tau)d\tau \leq v'(t) + k(t, t)v(t) \leq \\ &f(t)m(t) + k(t, t)v(t) \leq m(t)(f(t) + k(t, t)m(t)). \end{aligned} \quad (6)$$

Integrate (6) from  $a$  to  $t$  we obtain

$$m(t) \leq c \exp \left( \int_a^t (f(s) + k(s, s))ds \right). \quad (7)$$

Substitute (7) into (4) we get

$$v'(t) \leq cf(t) \exp \left( \int_a^t (f(s) + k(s, s))ds \right). \quad (8)$$

Integrating both sides of (8) from  $a$  to  $t$  we obtain

$$v(t) \leq c \left[ 1 + \int_a^t f(s) \exp \left( \int_a^s (f(\tau) + k(\tau, \tau))d\tau \right) ds \right]. \quad (9)$$

By (3) we have the desired result.

The assumption  $k_t(t, s) \in C[I \times I, R_-]$  which absents in [5] is necessary for validity of Theorem 2.1.

**Remark 2.2** If in Theorem 2.1. we put  $k(t, s) = g(s)$  we get the Pachpate's result [6] .

Now we give more precise and general version of Theorem 2.1.



**Theorem 2.3** Let  $u(t), f(t), a(t) \in C[R_+, R_+]$ ,  $k(t, s), k_t(t, s) \in C[G, R_+]$ ,  $G = \{(t, s) \in R_+^2 : 0 \leq s \leq t < \infty\}$  and  $c$  be a nonnegative constant.

If

$$u(t) \leq c + \int_0^t f(s)u(s)ds + \int_a^t f(s) \left( \int_0^s k(s, \tau)u(\tau)d\tau \right) ds, \quad (10)$$

for  $t \in R_+$  then

$$u(t) \leq c \left[ 1 + \int_0^t f(s) \exp \left( \int_0^s (f(\tau) + A(\tau))d\tau \right) ds \right], \quad (11)$$

for  $t \in R_+$  where

$$A(t) = k(t, t) + \int_0^t k_t(t, s)ds. \quad (12)$$

If

$$u(t) \leq a(t) + \int_0^t f(s)u(s)ds + \int_a^t f(s) \left( \int_0^s k(s, \tau)u(\tau)d\tau \right) ds, \quad (13)$$

for  $t \in R_+$  then

$$u(t) \leq a(t) + g(t) \left[ 1 + \int_0^t f(s) \exp \left( \int_0^s (f(\tau) + A(\tau))d\tau \right) ds \right], \quad (14)$$

for  $t \in R_+$  where

$$g(t) = \int_0^t f(s) \left[ a(s) + \int_0^s k(s, \tau)a(\tau)d\tau \right] ds, \quad (15)$$

$A(t)$  is defined by (12).

**Proof.** Define a function  $z(t)$  by the right hand side of (10). Then  $z(0) = c$ ,  $u(t) \leq z(t)$  and

$$z'(t) = f(t) \left[ u(t) + \int_0^t k(t, \tau)u(\tau)d\tau \right] \leq f(t) \left[ z(t) + \int_0^t k(t, \tau)z(\tau)d\tau \right]. \quad (16)$$

Define a function  $v(t)$  by

$$v(t) = z(t) + \int_0^t k(t, \tau)z(\tau)d\tau. \quad (17)$$

Then  $v(0) = z(0) = c$ ,  $z(t) \leq v(t)$ ,

$$z'(t) \leq f(t)v(t) \quad (18)$$

and  $v(t)$  is nondecreasing in  $t$  and

$$\begin{aligned} v'(t) &= z'(t) + k(t, t)z(t) + \int_0^t k_t(t, \tau)z(\tau)d\tau \\ &\leq f(t)v(t) + k(t, t)v(t) + \int_0^t k_t(t, \tau)z(\tau)d\tau \\ &\leq v(t) \left[ f(t) + k(t, t) + \int_0^t k_t(t, \tau)d\tau \right] \\ &= v(t)[f(t) + A(t)]. \end{aligned}$$

Thus

$$v(t) \leq c \exp \left( \int_0^s [f(\tau) + A(\tau)] d\tau \right). \quad (19)$$

Substituting (19) in (18) and integrating the resulting inequality from 0 to  $t$ ,  $t \in R_+$  we get

$$z(t) \leq c \left[ 1 + \int_0^t f(s) \exp \left( \int_0^s [f(\tau) + A(\tau)] d\tau \right) ds \right] \quad (20)$$

The desired inequality (11) follows from inequality  $u(t) \leq z(t)$ .

Now we prove inequality (14). Define a function  $w(t)$  by

$$w(t) = \int_0^t f(s) \left[ u(s) + \int_0^s k(s, \tau) u(\tau) d\tau \right] ds. \quad (21)$$

Then from (13)

$$u(t) \leq a(t) + w(t)$$

and using this in (21) we get

$$\begin{aligned} w(t) &\leq \int_0^t f(s) \left[ a(s) + w(s) + \int_0^s k(s, \tau) (a(\tau) + w(\tau)) d\tau \right] ds. \\ &= g(t) + \int_0^t f(s) \left[ w(s) + \int_0^s k(s, \tau) w(\tau) d\tau \right] ds, \end{aligned} \quad (22)$$

where  $g(t)$  is defined by (15). Clearly  $g(t)$  is nonnegative, continuous and nondecreasing in  $t$ . First, we assume  $g(t) > 0$  for  $t \in R_+$ . From (22) we have

$$\frac{w(t)}{g(t)} \leq 1 + \int_0^t f(s) \left[ \frac{w(s)}{g(s)} + \int_0^s k(s, \tau) \frac{w(\tau)}{g(\tau)} d\tau \right] ds.$$

Now an application of the inequality (10) we get

$$\frac{w(t)}{g(t)} \leq \left[ 1 + \int_0^t f(s) \exp \left( \int_0^s [f(\tau) + A(\tau)] d\tau \right) ds \right]. \quad (23)$$

The desired inequality (14) follows from (23) and the fact

$$u(t) \leq a(t) + w(t).$$

If  $g(t) \geq 0$  we carry out the above procedure with  $g(t) + \epsilon$  instead of  $g(t)$  where  $\epsilon > 0$  is an arbitrary small constant and then subsequently pass to the limit as  $\epsilon \rightarrow 0$  to obtain (14).

Now we give an application which is just sufficient to convey the importance of our results. Consider the nonlinear integrodifferential equation

$$x'(t) = f(t, x(t)) + \int_a^t g(t, s, x(s)) ds \quad (24)$$

and the corresponding perturbed equation

$$u'(t) = f(t, u(t)) + \int_a^t g(t, s, u(s))ds + h\left(t, u(t), \int_a^t k(t, s, u(s))ds\right), \quad (25)$$

where  $a, t \in R^+$ ,  $f \in C[R^+ \times R, R]$ ,  $g, k, h \in C[R^+ \times R^+ \times R, R]$ . Let  $x(t) = x(t, a, x_0)$  and  $u(t) = u(t, a, x_0)$  be solutions of (24) and (25) respectively with  $x(a) = u(a) = x_0$ . Suppose that  $f_x, g_x$  exist and  $f_x \in C[R^+ \times R, R]$ ,  $g_x \in C[R^+ \times R^+ \times R, R]$ . Put  $\Phi(t, a, x_0) = \frac{\partial x}{\partial x_0}(t, a, x_0)$  then (see[2])we get the following variational equations

$$z'(t) = f_x(t, x(t, a, x_0))z(t) + \int_a^t g_x(t, s, x(s, a, x_0))z(s)ds, \quad z(a) = 1 \quad (26)$$

$$\frac{\partial x}{\partial a}(t, a, x_0) + \Phi(t, a, x_0)f(a, x_0) \int_a^t \Phi(t, s, x_0)g(s, a, x_0)ds = 0 \quad (27)$$

and according to the nonlinear variation of constants formula we obtain

$$u(t) = x(t) + \int_a^t \Phi(t, s, u(s))h\left(s, u(s), \int_a^s k(s, \tau, u(\tau))d\tau\right)ds. \quad (28)$$

**Theorem 4.** Suppose that the following inequalities hold:

$$|\Phi(t, s, u)h(s, u, z)| \leq p(s)(|u| + |z|), \quad (29)$$

$$|k(t, s, u)| \leq q(t, s)|u| \quad (30)$$

for  $0 \leq s \leq t$ ,  $u, z \in R$ . If  $p(t)$  and  $q(t, s)$  are continuous nonnegative functions,  $q_t(t, s)$  is a nonpositive function and

$$\int_a^\infty p(s)ds < \infty, \quad \int_a^\infty q(s, s)ds < \infty \quad (31)$$

then for any bounded solution  $x(t, a, x_0)$  of (24) in  $R^+$  the corresponding solution  $u(t, a, x_0)$  of (25) is bounded in  $R^+$ .

**Proof.** As  $x(t, a, x_0)$  is bounded then  $|x(t)| \leq M$ , where  $M \in R^+$  is a nonzero constant. From (29),(30) and relation(28) we get

$$|u(t)| \leq M + \int_a^t p(s)|u(s)|ds + \int_a^t p(s) \left( \int_a^s q(s, \tau)|u(\tau)|d\tau \right) ds.$$

From Theorem 2.1 it follows

$$|u(t)| \leq M \left[ 1 + \int_a^t p(s) \exp \left( \int_a^s (p(\tau) + q(\tau, \tau))d\tau \right) ds \right].$$

According to (31) it follows that  $|u(t)|$  is bounded and the proof is complete.

## Acknowledgement

This research has been supported by the Czech Ministry of Education in the frames of MSM002160503 Research Intention MIKROSYN New Trends in Microelectronic Systems and Nanotechnologies and MSM0021630529 Research Intention Intelligent Systems in Automatization.

## References

- [1] BAINOV, D: SIMENEOV, P. *Integral inequalities and Applications*, Academic Publishers, Dordrecht, 1992.
- [2] BRAUER, F: *A nonlinear variation of constants formula for Volterra equations*, Mat. Systems Th., 6, 1972, 226-234.
- [3] CHANDRA, J., FLEISHMAN, B.A: *On a generalization of Gronwall-Bellman lemma in partially ordered Banach spaces*, J. Math. Appl., 31, 1970, 668-681.
- [4] OGUNTUASE, J.A: *Remarks on Gronwall type inequalities*, An. Stiint. Univ. "Al. I. Cuza", 45, 1999, 321-328.
- [5] OGUNTUASE, J.A: *On an inequality of Gronwall*, J. of Inequalities in Pure and Applied Math., 2, 2001, 1-6.
- [6] PACHPATTE, B.G: *A note on Gronwall-Bellman inequality*, J. Math. Anal. Appl., 44, 1973, 758-762.
- [7] PACHPATTE, B.G: *Inequalities for Differential and Integral Equations*, Mathematics in Science and Engineering, vol. 197, Academic Press Inc, San Diego, 2006.
- [8] PACHPATTE, B.G: *Mathematical Inequalities*, Elsevier, Amsterdam, vol.67, 2005.
- [9] ŠMARDÁ, Z.: *Modifications of the Gronwall inequality and their application*, XXVI. Proceedings of International Coloquium , Brno, 2008, 1-5.

## Current address

**Doc. RNDr. Zdeněk Šmarda, CSc.**

Department of Mathematics

Faculty of Electrical Engineering and Communication

Brno University of Technology, Technická 8, 616 00 BRNO

e-mail: smarda@feec.vutbr.cz

# CONVERGENCE PROOF OF A MONTE CARLO SCHEME FOR THE RESOLUTION OF THE SMOLUCHOWSKI COAGULATION EQUATION

BARAKEH Bilal, (RL)

**Abstract.** This paper studies the convergence in probability of a Monte Carlo simulation scheme for solving Smoluchowski's coagulation equation. We propose the lemmas and theorems needed to achieve the convergence proof.

**Key words:** Smoluchowski equation ; Monte Carlo scheme ; Convergence in probability.

## 1 Introduction

The mathematical foundation and numerical simulation of cluster growth is a research field of high current interest. Models of cluster growth arise in various phenomena and find their applications in a wide range of engineering contexts ranging from environmental sciences (growing and spreading of air pollutants) to the development of engines (behaviour of fuel mixtures). This theory models and describes the evolution of a system of a large number of particles that can coagulate to form clusters which in turn can coalesce in order to form bigger clusters. Each cluster is identified by its size. In his work on coagulation processes in colloids, M.V. Smoluchowski [1] proposed an infinite system of differential equations to describe the time evolution in some physical system of the concentration of particles of size  $i$  at time  $t$  denoted by  $c(i, t)$ . In the simplest situation, the space homogeneous discrete Smoluchowski equation reads, for  $i=1,2,3,\dots$  and  $t>0$  :

$$\begin{cases} \frac{\partial c}{\partial t}(i, t) = \frac{1}{2} \sum_{j=1}^{i-1} K(i-j, j) c(i-j, t) c(j, t) - \sum_{j=1}^{\infty} K(i, j) c(i, t) c(j, t) \\ c(i, 0) = c_i(0) \end{cases} \quad (1)$$

In fact, this system describes a nonlinear evolution equation of infinite dimension, with initial condition  $(c_i(0))_{i \geq 1}$  and so that  $\sum_{i \geq 1} c_i(0) = 1$ . The rate of merging of particles of size  $i$  and  $j$  at time  $t$  is given by the coagulation kernel  $K(i, j)$  that is naturally supposed to be nonnegative

[i.e.,  $K : (IN^*)^2 \rightarrow IR^+$ ] and symmetric [i.e.,  $K(i, j) = K(j, i)$ ]. The structure of equation (1) is closely related to that of the fundamental Boltzmann equation of rarefied gas dynamics with a transport differential operator on the left hand side and a local quadratic particle interaction operator on the right hand side. The first term on the right hand side shows that the concentration of particles of size  $i$  increases as a result of coagulation of particles of sizes  $i-j$  and  $j$ . This is the gain term. The coefficient  $\frac{1}{2}$  is due to the fact that  $K$  is symmetric. The second term corresponds to the depletion of particles of size  $i$  after coalescence with other particles. It represents the loss term. Due to the complexity of the Smoluchowski equation, the problem of numerically solving it with deterministic methods is a difficult task that cannot be handled on a computer with reasonable calculation efforts. The main numerical tools for solving such equations are Monte Carlo simulations. Several stochastic algorithms have been proposed [2,3,4] and are now understood as mathematically rigorous numerical algorithms. But the convergence of the simulated solutions is a field where researchers are still contributing to improve it [5,6]. The aim of this work is to study the convergence of a modified Monte Carlo scheme proposed in [2]. We first reformulate the algorithm and then we propose the lemmas and theorems needed to achieve the convergence proof.

## 2 Reformulation of the algorithm

The algorithm proposed in [2] is based on a Monte Carlo scheme that takes a system of test particles which interact and form clusters according to the dynamics described above. Random numbers are used to find out which clusters interact and to determine the size of the new clusters. We first recall some basic notations and concepts. If  $N_0$  is the initial total number of particles, then at time  $t$ ,  $N_0 c(i, t)$  represents the total number of particles of size  $i$  and  $\sum_{i \geq 1} N_0 c(i, t)$  is the total number of particles. The quantity  $ic(i, t)$  represents the fraction at time  $t$  of the whole mass produced by the particles of mass  $i$ . The whole mass is given by

$$m(t) = \sum_{i \geq 1} ic(i, t).$$

Multiplying equation (1) by  $i$  and summing over all, one can verify that the whole mass is conserved

$$\frac{d}{dt} \sum_{i \geq 1} ic(i, t) = 0, \quad (2)$$

provided that the relevant summations converge and can be interchanged which is valid as long as

$$\sum_{i, j \geq 1} K(i, j) c(i, t) c(j, t) < \infty$$

We refer to [7] for a study of existence, uniqueness, and conservation of mass of solutions. Since particles may stick together, the total number of particles is a decreasing quantity and this may display a poor statistics for the simulation. Therefore, to avoid this deficiency one can approximate the following mass density function

$$g(i, t) = ic(i, t), \quad (3)$$

If we write  $\tilde{g}_i(t)$  instead of  $g(i, t)$  the equation (1) becomes:

$$\begin{cases} \frac{\partial \tilde{g}_i}{\partial t}(t) = \sum_{j=1}^{i-1} K^*(i-j, j) \tilde{g}_{i-j}(t) \tilde{g}_j(t) - \sum_{j=1}^{\infty} K^*(i, j) \tilde{g}_i(t) \tilde{g}_j(t) \\ \tilde{g}_i(0) = g_{0,i} \end{cases} \quad (4)$$

Here,

$$K^*(i, j) = \frac{K(i, j)}{j}.$$

The kernel  $K^*$  is bounded and without loss of generality we assume that at time  $t=0$

$$\sum_{i \geq 1} g_{0,i} = 1. \quad (5)$$

Therefore, the conservation of mass leads to

$$\sum_{i \geq 1} \tilde{g}_i(t) = 1. \quad (6)$$

We start with an initial  $N$ -tuple  $Z_0^{(N)} = (z_{0,1}^{(N)}, z_{0,2}^{(N)}, \dots, z_{0,N}^{(N)}) \in IN^N$ . The entry  $z_{0,i}^{(N)}$  is representing a particle of mass  $1/N$  with 'label'  $i$  and such that

$$\forall i \in IN^* \quad \frac{1}{N} \# \{ \ell : z_{0,\ell}^{(N)} = i \} \approx \tilde{g}_i(0) \quad (7)$$

If we assume a monodisperse initial condition

$$\tilde{g}_1(0) = 1, \tilde{g}_2(0) = \tilde{g}_3(0) = \dots = 0, \quad (8)$$

we set,

$$z_{0,1}^{(N)} = z_{0,2}^{(N)} = \dots = z_{0,N}^{(N)} = 1. \quad (9)$$

We choose a time step  $\Delta t$  such that

$$\Delta t \sup_{i,j \in IN} K^*(i, j) < 1. \quad (10)$$

This means that the time step has to be permanently reduced while particles progress in to the domain of increasing mass. For  $n \in IN$ , we set  $t_n = n\Delta t$ . At  $t = t_n$ , we consider a point set  $Z^{(N)}(n)$  of  $N$  particles  $z_1^{(N)}(n), \dots, z_N^{(N)}(n)$  such that  $\forall i \in IN^*$ ,

$$\frac{1}{N} \# \{ \ell : z_\ell^{(N)}(n) = i \} \approx \tilde{g}_i(t_n). \quad (11)$$

For  $i \in IN^*$ , we define the following independent equally distributed random numbers

$$\chi_{i,\ell}^{(N)}(n) = \begin{cases} 1 & \text{si } z_\ell^{(N)}(n) = i, \\ 0 & \text{si non.} \end{cases} \quad (12)$$

such that

$$G_i^{(N)}(n) = \frac{1}{N} \sum_{\ell=1}^N \chi_{i,\ell}^{(N)}(n), \quad (13)$$

with,

$$\sum_{i=1}^{\infty} G_i^{(N)}(n) = 1. \quad (14)$$

One can notice easily that

$$G_i^{(N)}(n) = \frac{1}{N} \# \{ \ell : z_{\ell}^{(N)}(n) = i \} \quad (15)$$

Applying the Euler time discretization for  $\tilde{g}_i(t_n)$  will lead to

$$\begin{aligned} \frac{1}{\Delta t} (\tilde{g}_i(t_{n+1}) - \tilde{g}_i(t_n)) &= \sum_{j=1}^{i-1} K^*(i-j, j) \tilde{g}_{i-j}(t_n) \tilde{g}_j(t_n) \\ &\quad - \sum_{j=1}^{\infty} K^*(i, j) \tilde{g}_i(t_n) \tilde{g}_j(t_n), \quad i \in \mathbb{N}^*, n \in \mathbb{N}. \end{aligned}$$

Using now the fact that,  $\forall n \in \mathbb{N}$

$$\sum_{j \geq 1} \tilde{g}_j(t_n) = 1,$$

we conclude that  $\tilde{g}_i(t_{n+1})$  is defined by

$$\tilde{g}_i(t_{n+1}) = \sum_{j=1}^{i-1} \Delta t K^*(i-j, j) \tilde{g}_{i-j}(t_n) \tilde{g}_j(t_n) + \sum_{j=1}^{\infty} (1 - \Delta t K^*(i, j)) \tilde{g}_i(t_n) \tilde{g}_j(t_n) \quad (16)$$

Then the Monte Carlo scheme proposed in [2] can now be announced as follows:

**Initialisation step:** For  $1 \leq i \leq N$ , choose at time  $t=0$ ,  $z_i^{(N)}(0) \in \{1, 2, 3, \dots\}$  such that

$$\forall i \in \mathbb{N}^* \quad G_i^{(N)}(0) \approx g_{0,i}$$

**Coagulation step:** For  $1 \leq \ell \leq N$ , the transition from  $z_{\ell}^{(N)}(n)$  à  $z_{\ell}^{(N)}(n+1)$  is given by the following random game:

(i) For  $i=1, \dots, N$  choose equidistributed random numbers  $\pi_i^{(N)} \in \{1, 2, \dots, N\}$  and  $r_i^{(N)} \in [0, 1]$ .

(ii) Choose a time step  $\Delta t$  such that  $\Delta t \sup_{i,j \in \mathbb{N}} K^*(i, j) < 1$ .

(iii) Let  $p(i, j) = \Delta t K^*(i, j)$  and define for  $1 \leq \ell \leq N$

$$z_{\ell}^{(N)}(n+1) := \begin{cases} z_{\ell}^{(N)}(n) + z_{\pi_{\ell}^{(N)}(n)}^{(N)}(n) & \text{if } r_{\ell}^{(N)}(n) \leq p(z_{\ell}^{(N)}(n), z_{\pi_{\ell}^{(N)}(n)}^{(N)}(n)), \\ z_{\ell}^{(N)}(n) & \text{if not} \end{cases} \quad (17)$$



### 3 Convergence proof

To establish the convergence of the numerical scheme, we use the following two lemmas:

**Lemma 1.** Let  $(\Omega, A, P)$  be a probability space. Let  $(X_N)_{N \geq 1}$  be a sequence of real random variables into  $L^2(\Omega)$ . If  $E[X_N]$  is the expected value of the variable  $X_N$  and  $\gamma \in \mathbb{R}$  then the following two conditions are equivalent :

$$(a) X_N \xrightarrow{L^2(\Omega)} \gamma \quad ; \quad (b) \lim_{N \rightarrow \infty} E[X_N] = \gamma \text{ and } \lim_{N \rightarrow \infty} E[X_N^2] = \gamma^2$$

**Lemma 2.** Let  $(P_N)_{N \geq 1}$  be a sequence of probabilities on  $IN^*$ , defined by

$$P_N = \sum_{i=1}^{\infty} \alpha_i^{(N)} \delta^{\{i\}},$$

which converges weakly<sup>(1)</sup> to the probability  $P$  on  $IN^*$ , with

$$P = \sum_{i=1}^{\infty} \alpha_i \delta^{\{i\}},$$

where,  $\delta^{\{i\}}$  is the Dirac measure on  $IN^*$  and  $\sum_{i=1}^{\infty} \alpha_i^{(N)} = 1$ .

If  $(s^{(N)})_{N \geq 1}$  is a family of uniformly bounded sequences and if  $(s)$  is a bounded sequence, such that

$$\forall i \in IN^* \quad \lim_{N \rightarrow \infty} s^{(N)}(i) = s(i),$$

then

$$\lim_{N \rightarrow \infty} \sum_{i=1}^{\infty} \alpha_i^{(N)} s^{(N)}(i) = \sum_{i=1}^{\infty} \alpha_i s(i).$$

In particular, this implies that for all bounded sequence  $\sigma$ , we have

$$\lim_{N \rightarrow \infty} \sum_{i=1}^{\infty} \alpha_i^{(N)} \sigma(i) = \sum_{i=1}^{\infty} \alpha_i \sigma(i),$$

The convergence of the numerical scheme is resulting from the following proposition that will be announced and verified.

**Proposition.** For all  $n \in IN$ , if

$$\forall i \in IN^* : G_i^{(N)}(n) \xrightarrow{L^2} \tilde{g}_i(t_n),$$

then

$$\forall i \in IN^* : G_i^{(N)}(n+1) \xrightarrow{L^2} \tilde{g}_i(t_{n+1}).$$

**Proof:**

By the lemma 1, we notice that we just need to show that : For  $n \in IN$  and  $i \in IN^*$ ,

$$\lim_{N \rightarrow \infty} E[G_i^{(N)}(n+1)] = \tilde{g}_i(t_{n+1}) \quad (18)$$

---

<sup>(1)</sup>  $P_N \xrightarrow[N \rightarrow \infty]{\text{Weakly}} P$  if,  $\forall i \in IN^* : \lim_{N \rightarrow \infty} \alpha_i^{(N)} = \alpha_i$  (i.e.)  $\forall \varphi \in \mathcal{C}_0 : \langle P_N, \varphi \rangle \xrightarrow[N \rightarrow \infty]{} \langle P, \varphi \rangle$

and

$$\lim_{N \rightarrow \infty} E[\{G_i^{(N)}(n+1)\}^2] = \{\lim_{N \rightarrow \infty} E[G_i^{(N)}(n+1)]\}^2 = \{\tilde{g}_i(t_{n+1})\}^2 \quad (19)$$

- To verify (18), we use the formula of conditional expectations. For  $1 \leq \ell \leq N$ , we have

$$E[\chi_{i,\ell}^{(N)}(n+1)] = P(z_\ell^{(N)}(n+1) = i) = \sum_{j=1}^{\infty} \sum_{k=1}^{\infty} \underbrace{P(z_\ell^{(N)}(n+1) = i \mid z_\ell^{(N)}(n) = j, z_{\pi_\ell^{(N)}(n)}^{(N)}(n) = k)}_{=p_1} \times \underbrace{P(z_\ell^{(N)}(n) = j, z_{\pi_\ell^{(N)}(n)}^{(N)}(n) = k)}_{=p_2}.$$

We compute now the terms  $p_1$  and  $p_2$  as follows:

**a)** if  $j < i$  and  $k = i-j$ , then

$$\begin{aligned} p_1 &= P(z_\ell^{(N)}(n+1) = i \mid z_\ell^{(N)}(n) = j, z_{\pi_\ell^{(N)}(n)}^{(N)}(n) = k) \\ &= P(r_\ell^{(N)}(n) \leq p(j, i-j) \mid z_\ell^{(N)}(n) = j, z_{\pi_\ell^{(N)}(n)}^{(N)}(n) = i-j) \\ &= p(j, i-j); \end{aligned}$$

**b)** if  $j = i$ , then

$$\begin{aligned} p_1 &= P(z_\ell^{(N)}(n+1) = i \mid z_\ell^{(N)}(n) = j, z_{\pi_\ell^{(N)}(n)}^{(N)}(n) = k) \\ &= P(z_\ell^{(N)}(n+1) = z_\ell^{(N)}(n) \mid z_\ell^{(N)}(n) = i, z_{\pi_\ell^{(N)}(n)}^{(N)}(n) = k) \\ &= P(r_\ell^{(N)}(n) > p(i, k) \mid z_\ell^{(N)}(n) = i, z_{\pi_\ell^{(N)}(n)}^{(N)}(n) = k) \\ &= 1 - p(i, k); \end{aligned}$$

**c)** otherwise,

$$p_1 = P(z_\ell^{(N)}(n+1) = i \mid z_\ell^{(N)}(n) = j, z_{\pi_\ell^{(N)}(n)}^{(N)}(n) = k) = 0.$$

As for the term  $p_2$ , we have

$$\begin{aligned} p_2 &= P(z_\ell^{(N)}(n) = j, z_{\pi_\ell^{(N)}(n)}^{(N)}(n) = k) \\ &= \sum_{m=1}^N P(z_\ell^{(N)}(n) = j, z_m^{(N)}(n) = k, \pi_\ell^{(N)}(n) = m) = \sum_{m=1}^N P(z_\ell^{(N)}(n) = j) P(z_m^{(N)}(n) = k) P(\pi_\ell^{(N)}(n) = m) \\ &= P(z_\ell^{(N)}(n) = j) P(\pi_\ell^{(N)}(n) = m) \sum_{m=1}^N P(z_m^{(N)}(n) = k) \\ &= E[\chi_{j,\ell}^{(N)}(n)] \frac{1}{N} \sum_{m=1}^N E[\chi_{k,m}^{(N)}(n)] \\ &= E[\chi_{j,\ell}^{(N)}(n)] E[G_k^{(N)}(n)] \end{aligned}$$

Therefore, for  $1 \leq \ell \leq N$

$$E[\chi_{i,\ell}^{(N)}(n+1)] = \sum_{j=1}^{i-1} p(j, i-j) E[\chi_{j,\ell}^{(N)}(n)] E[G_{i-j}^{(N)}(n)] + E[\chi_{i,\ell}^{(N)}(n)] \sum_{k=1}^{\infty} (1 - p(i, k)) E[G_k^{(N)}(n)].$$

If we divide the last equality by  $N$  and we sum over all  $\ell = 1$  to  $N$ , we obtain

$$E[G_i^{(N)}(n+1)] = \sum_{j=1}^{i-1} p(j, i-j) E[G_j^{(N)}(n)] E[G_{i-j}^{(N)}(n)] + E[G_i^{(N)}(n)] \sum_{k=1}^{\infty} (1-p(i, k)) E[G_k^{(N)}(n)]$$

With the relations (6) and (14), let us define on  $IN^*$  the probabilities:

$$P = \sum_{i=1}^{\infty} \tilde{g}_j(t_n) \delta^{\{i\}} \text{ and } P_N = \sum_{i=1}^{\infty} E[G_j^{(N)}(n)] \delta^{\{i\}}.$$

Using the assumption of the proposition and the lemma, we deduce that

$$\forall j \in IN^* \quad \lim_{N \rightarrow \infty} E[G_j^{(N)}(n)] = \tilde{g}_j(t_n), \quad (20)$$

hence, the sequence  $(P_N)_{N \geq 1}$  converges in the weak sense to  $P$ . Now, if from one side we define the sequence

$$s(j) = \begin{cases} p(j, i-j) \tilde{g}_{i-j}(t_n) & \text{if } j < i, \\ 0 & \text{if not} \end{cases} \quad (21)$$

and from another side, we introduce for  $N \in IN^*$  the uniformly bounded sequences :

$$s^{(N)}(j) = \begin{cases} p(j, i-j) E[G_{i-j}^{(N)}(n)] & \text{if } j < i, \\ 0 & \text{if not} \end{cases} \quad (22)$$

We deduce from (20) and lemma (2) that

$$\lim_{N \rightarrow \infty} \sum_{j=1}^{i-1} E[G_j^{(N)}(n)] p(j, i-j) E[G_{i-j}^{(N)}(n)] = \sum_{j=1}^{i-1} \tilde{g}_j(t_n) p(j, i-j) \tilde{g}_{i-j}(t_n)$$

Using now the sequence  $\sigma(k) = 1 - p(i, k)$ , one can verify also with (20) and lemma (2) that

$$\lim_{N \rightarrow \infty} \sum_{k=1}^{\infty} E[G_k^{(N)}(n)] (1 - p(i, k)) = \sum_{k=1}^{\infty} \tilde{g}_k(t_n) (1 - p(i, k)).$$

This allows us in particular to write

$$\lim_{N \rightarrow \infty} E[G_i^{(N)}(n+1)] = \sum_{j=1}^{i-1} \tilde{g}_j(t_n) p(j, i-j) \tilde{g}_{i-j}(t_n) + \tilde{g}_i(t_n) \sum_{j=1}^{i-1} \tilde{g}_k(t_n) (1 - p(i, k)).$$

By remarking that the right hand side of this equation is equal to  $\tilde{g}_i(t_{n+1})$ , we finally conclude that  $\lim_{N \rightarrow \infty} E[G_i^{(N)}(n+1)] = \tilde{g}_i(t_{n+1})$ . This will end the proof of (18).

- As for the relation (19), we will prove it as follows:

$$\{G_i^{(N)}(n+1)\}^2 = \frac{1}{N^2} \sum_{\ell=1}^N \chi_{i,\ell}^{(N)}(n+1) + \frac{1}{N^2} \sum_{\substack{\ell, m=1 \\ \ell \neq m}}^N \chi_{i,\ell}^{(N)}(n+1) \chi_{i,m}^{(N)}(n+1),$$

(the square of the indicator function is equal to itself).

But for  $\ell \neq m$ , the variables  $\chi_{i,\ell}^{(N)}(n+1)$  and  $\chi_{i,m}^{(N)}(n+1)$  are independent. Hence,

$$E[\{G_i^{(N)}(n+1)\}^2] = \frac{1}{N^2} \sum_{\ell=1}^N E[\chi_{i,\ell}^{(N)}(n+1)] + \frac{1}{N^2} \sum_{\substack{\ell,m=1 \\ \ell \neq m}}^N E[\chi_{i,\ell}^{(N)}(n+1)]E[\chi_{i,m}^{(N)}(n+1)],$$

Since the variables  $\chi_{i,\ell}^{(N)}(n+1)$  are Bernoulli's random variables, one can deduce easily that :

$$\frac{1}{N^2} \sum_{\ell=1}^N E[\chi_{i,\ell}^{(N)}(n+1)] \leq \frac{1}{N} \quad \text{and} \quad \frac{1}{N^2} \sum_{\ell=1}^N \{E[\chi_{i,\ell}^{(N)}(n+1)]\}^2 \leq \frac{1}{N}$$

hence,

$$\lim_{N \rightarrow \infty} \frac{1}{N^2} \sum_{\ell=1}^N E[\chi_{i,\ell}^{(N)}(n+1)] = 0 = \lim_{N \rightarrow \infty} \frac{1}{N^2} \sum_{\ell=1}^N \{E[\chi_{i,\ell}^{(N)}(n+1)]\}^2$$

This means,

$$\begin{aligned} \lim_{N \rightarrow \infty} E[(G_i^{(N)}(n+1))^2] &= \lim_{N \rightarrow \infty} \left\{ \frac{1}{N^2} \sum_{\ell=1}^N (E[\chi_{i,\ell}^{(N)}(n+1)])^2 + \frac{1}{N^2} \sum_{\substack{\ell,m=1 \\ \ell \neq m}}^N E[\chi_{i,\ell}^{(N)}(n+1)]E[\chi_{i,m}^{(N)}(n+1)] \right\} \\ &= \lim_{N \rightarrow \infty} \left\{ \frac{1}{N} \sum_{\ell=1}^N E[\chi_{i,\ell}^{(N)}(n+1)] \right\}^2, \end{aligned}$$

which concludes the proof of (19).

Finally, one can see that the proof of the proposition is resulting from the lemma (1) and concludes the convergence proof of the algorithm.

## 4 Conclusions

In this paper we analysed a procedure for solving Smoluchowski's coagulation equation. A detailed convergence proof based on a probabilistic approach has been formulated. This work aims to be a contribution in the study of the convergence of other similar numerical schemes.

## References

- [1.] SMOLUCHOWSKI, M.V. 1916 Versuch einer mathematischen Theorie der Koagulationskinetik kolloider Lösungen. *Z. Phys.Chem.* **92**, 129-168.
- [2.] BABOVSKY, H. 1999 On a Monte Carlo scheme for Smoluchowski's coagulation equation. *Monte Carlo Methods Appl.* **5**, 1-18.
- [3.] EIBECK, A & WAGNER, W. 2000 An efficient stochastic algorithm for studying coagulation dynamics and gelation phenomena. *Siam J.Sci.Comput.* **22**, 802-821.
- [4.] SABELFELD, K.K, ROGASINSKY, S.V, KOLODKA, A.A and Levykin A.I. (1996), Stochastic algorithms for solving Smoluchowski coagulation equation & applications to aerosol growth simulation. *Monte Carlo Methods Appl.* **2**, 41-87.
- [5.] EIBECK, A & WAGNER, W. 2001 Stochastic particle approximations for Smoluchowski's coagulation equation. *The Annals of Applied Probability.* **11**, 1137-1165.
- [6.] LÉCOT, C. & WAGNER, W. 2004 A quasi Monte Carlo scheme for Smoluchowski's coagulation equation. *Mathematics of Computation.* **73**, 1953-1966.

- [7.] HEILMANN, O.J. 1992 Analytical solutions of Smoluchowski's coagulation equation.  
*J.Phys. A* **25**, 3763-3771.

**Current address**

**BARAKEH Bilal**

Assistant professor, University of Balamand, Faculty of Sciences,  
Department of Mathematics. B.P:100-Tripoli,Lebanon,  
e-mail : bilal.barake@balamand.edu.lb



# RECONSTRUCTION OF CLOSELY SPACED SMALL INHOMOGENEITIES VIA BOUNDARY MEASUREMENTS FOR THE FULL TIME-DEPENDENT MAXWELL'S EQUATIONS

DAVEAU Christian, (F), KHELIFI Abdessatar, (F), SUSHCHENKO Anton, (F)

**Abstract.** We consider for the full time-dependent Maxwell's equations the inverse problem of identifying locations and certain properties of small electromagnetic inhomogeneities in a homogeneous background medium from dynamic boundary measurements on the boundary for a finite time interval.

**Key words and phrases.** Maxwell's equations, inhomogeneities, inverse problem, reconstruction, geometric control

*Mathematics Subject Classification.* 35R30, 35B40, 35B37, 78M35

## 1 Introduction

The ultimate objective of the work described in this paper is to determine locations and certain properties of the shapes of small electromagnetic inhomogeneities in a homogeneous background medium from dynamic boundary measurements on part of the boundary and for finite interval in time. Using as weights particular background solutions constructed by a geometrical control method we develop an asymptotic method based on appropriate averaging of the partial dynamic boundary measurements.

For stationary Maxwell's equations it has been known that the Dirichlet to Neumann map uniquely determines (smooth) isotropic electromagnetic parameters, see [16], [18], [20]. We will provide in this paper a rigorous derivation of the inverse Fourier transform of a linear combination of derivatives of point masses, located at the positions  $z_j$  of the inhomogeneities, as the

leading order term of an appropriate averaging of (partial) dynamic boundary measurements of the tangential components of electric fields on part of the boundary. Our formulas may be used to determine properties (location, relative size ) of the small inhomogeneities in case a single, or a few (tangential) boundary electric fields are known. Our approach differs from [1], [2], [3], [4], [22] and is expected to lead to very effective computational identification algorithms. Our main result is given by:

**Theorem 4.1** *Let  $\eta \in \mathbb{R}^d$ . Let  $E_\alpha$  be the unique solution in  $\mathcal{C}^0(0, T; X(\Omega)) \cap \mathcal{C}^1(0, T; L^2(\Omega))$  to the Maxwell's equations (3) with  $\varphi(x) = \eta^\perp e^{i\eta \cdot x}$ ,  $\psi(x) = -i\sqrt{\mu_0}|\eta|\eta^\perp e^{i\eta \cdot x}$ , and  $f(x, t) = \eta^\perp e^{i\eta \cdot x - i\sqrt{\mu_0}|\eta|t}$ . Suppose that  $\Gamma$  and  $T$  geometrically control  $\Omega$ , then we have*

$$\int_0^T \int_\Gamma \left[ \theta_\eta \cdot (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}) + \partial_t \theta_\eta \cdot \partial_t (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}) \right] d\sigma(x) dt = \alpha^2 \sum_{j=1}^m (\mu_0 - \mu_j) e^{2i\eta \cdot z_j} M_j(\eta) \cdot \eta + O(\alpha^2),$$

where  $\theta_\eta$  is the unique solution to the Volterra equation (20) with  $g_\eta$  defined as the boundary control in (18) and  $M_j$  is the polarization tensor of  $B_j$ , defined by

$$(M_j)_{k,l} = e_k \cdot \left( \int_{\partial B_j} \left( \nu_j + \left( \frac{\mu_j}{\mu_0} - 1 \right) \frac{\partial \Phi_j}{\partial \nu_j} \Big|_+(y) \right) y \cdot e_l ds_j(y) \right).$$

Here  $(e_1, e_2)$  is an orthonormal basis of  $\mathbb{R}^d$ . The term  $O(\alpha^2)$  is independent of the points  $\{z_j, j = 1, \dots, m\}$ .

For discussions on closely related (stationary) identification problems we refer the reader to [19],[21], [6], and [10].

## 2 Problem formulation

Let  $\Omega$  be a bounded  $C^2$ -domain in  $\mathbb{R}^d$ ,  $d = 2, 3$ . Assume that  $\Omega$  contains a finite number of inhomogeneities, each of the form  $z_j + \alpha B_j$ , where  $B_j \subset \mathbb{R}^d$  is a bounded, smooth domain containing the origin. The total collection of inhomogeneities is  $\mathcal{B}_\alpha = \cup_{j=1}^m (z_j + \alpha B_j)$ . The points  $z_j \in \Omega$ ,  $j = 1, \dots, m$ , which determine the location of the inhomogeneities, are assumed to satisfy the following inequalities:

$$|z_j - z_l| \geq c_0 > 0, \forall j \neq l \quad \text{and} \quad \operatorname{dist}(z_j, \partial\Omega) \geq c_0 > 0, \forall j. \quad (1)$$

Assume that  $\alpha > 0$ , the common order of magnitude of the diameters of the inhomogeneities, is sufficiently small, that these inhomogeneities are disjoint, and that their distance to  $\mathbb{R}^d \setminus \overline{\Omega}$  is larger than  $c_0/2$ . Let  $\mu_0$  and  $\varepsilon_0$  denote the permeability and the permittivity of the background medium, and assume that  $\mu_0 > 0$  and  $\varepsilon_0 > 0$  are positive constants. Let  $\mu_j > 0$  and  $\varepsilon_j > 0$  denote the permeability and the permittivity of the  $j$ -th inhomogeneity,  $z_j + \alpha B_j$ , these are also assumed to be positive constants. Introduce the piecewise-constant magnetic permeability

$$\mu_\alpha(x) = \begin{cases} \mu_0, & x \in \Omega \setminus \overline{\mathcal{B}_\alpha}, \\ \mu_j, & x \in z_j + \alpha B_j, j = 1 \dots m. \end{cases} \quad (2)$$



If we allow the degenerate case  $\alpha = 0$ , then the function  $\mu_0(x)$  equals the constant  $\mu_0$ . The electric permittivity is defined by  $\varepsilon_\alpha(x) = \varepsilon_0$ , for all  $x \in \Omega$ . Let  $\mathbf{n} = \mathbf{n}(x)$  denote the outward unit normal vector to  $\Omega$  at a point on  $\partial\Omega$ ,  $\partial_t u = \frac{\partial u}{\partial t}$  and  $\Delta$  means the Laplace operator defined

$$\text{by } \Delta u = \sum_{i=1}^d \frac{\partial^2 u}{\partial x_i^2}.$$

In this paper, we will denote by bold letters the functional spaces for the vector fields. Thus  $H^s(\Omega)$  denotes the usual Sobolev space on  $\Omega$  and  $\mathbf{H}^s(\Omega)$  denotes  $(H^s(\Omega))^d$  and  $\mathbf{L}^2(\Omega)$  denotes  $(L^2(\Omega))^d$ . As usual for Maxwell equations, we need spaces of fields with square integrable curls:

$$\mathbf{H}(\text{curl}; \Omega) = \{u \in \mathbf{L}^2(\Omega), \text{curl } u \in \mathbf{L}^2(\Omega)\},$$

and with square integrable divergences

$$\mathbf{H}(\text{div}; \Omega) = \{u \in \mathbf{L}^2(\Omega), \text{div } u \in L^2(\Omega)\}.$$

We will also need the following functional spaces:

$$Y(\Omega) = \{u \in \mathbf{L}^2(\Omega), \text{div } u = 0 \text{ in } \Omega\}, \quad X(\Omega) = \mathbf{H}^1(\Omega) \cap Y(\Omega),$$

and  $TL^2(\partial\Omega)$  the space of vector fields on  $\partial\Omega$  that lie in  $\mathbf{L}^2(\partial\Omega)$ . Finally, the "minimal" choice for the electric variational space would be

$$X_N(\Omega) = \{v \in \mathbf{H}(\text{curl}; \Omega) \cap \mathbf{H}(\text{div}; \Omega); \quad v \times \mathbf{n} = 0 \quad \text{on } \partial\Omega\}.$$

Now, we introduce the following time-dependent Maxwell equations (associated to the electric field)

$$\begin{cases} (\varepsilon_\alpha \partial_t^2 + \text{curl } \frac{1}{\mu_\alpha} \text{curl}) E_\alpha = 0 & \text{in } \Omega \times (0, T), \\ \text{div } (\varepsilon_\alpha E_\alpha) = 0 & \text{in } \Omega \times (0, T), \\ E_\alpha|_{t=0} = \varphi, \partial_t E_\alpha|_{t=0} = \psi & \text{in } \Omega, \\ E_\alpha \times \mathbf{n}|_{\partial\Omega \times (0, T)} = f, \end{cases} \quad (3)$$

where  $E_\alpha \in \mathbb{R}^d$  is the electric field,  $f$  the boundary condition for  $E_\alpha \times \mathbf{n}$ , and  $\varphi$  and  $\psi$  the initial data.

Let  $E$  be the solution of the Maxwell's equations in the homogeneous domain:

$$\begin{cases} (\varepsilon_0 \partial_t^2 + \text{curl } \frac{1}{\mu_0} \text{curl}) E = 0 & \text{in } \Omega \times (0, T), \\ \text{div } (\varepsilon_0 E) = 0 & \text{in } \Omega \times (0, T), \\ E|_{t=0} = \varphi, \partial_t E|_{t=0} = \psi & \text{in } \Omega, \\ E \times \mathbf{n}|_{\partial\Omega \times (0, T)} = f. \end{cases} \quad (4)$$

Here  $T > 0$  is a final observation time and  $\varphi, \psi \in \mathcal{C}^\infty(\overline{\Omega})$  and  $f \in \mathcal{C}^\infty(0, T; \mathcal{C}^\infty(\partial\Omega))$  are subject to the compatibility conditions

$$\partial_t^{2l} f|_{t=0} = (\Delta^l \varphi) \times \mathbf{n}|_{\partial\Omega} \text{ and } \partial_t^{2l+1} f|_{t=0} = (\Delta^l \psi) \times \mathbf{n}|_{\partial\Omega}, \quad l = 1, 2, \dots$$

it follows that (4) has a unique solution  $E \in \mathcal{C}^\infty([0, T] \times \overline{\Omega})$ . It is also known (see for example [17]) that since  $\Omega$  is smooth ( $\mathcal{C}^2$ -regularity would be sufficient) the non homogeneous Maxwell's equations (3) have a unique weak solution  $E_\alpha \in \mathcal{C}^0(0, T; X(\Omega)) \cap \mathcal{C}^1(0, T; \mathbf{L}^2(\Omega))$ . Indeed,  $\text{curl } E_\alpha$  belongs to  $\mathcal{C}^0(0, T; X(\Omega)) \cap \mathcal{C}^1(0, T; \mathbf{L}^2(\Omega))$ .

### 3 Asymptotic formula

We start the derivation of the asymptotic formula for  $\text{curl } E_\alpha \times \mathbf{n}$  with the following estimate.

**Lemma 3.1** *The following estimate as  $\alpha \rightarrow 0$  holds:*

$$\|\partial_t(E_\alpha - E)\|_{L^\infty(0, T; \mathbf{L}^2(\Omega))} + \|E_\alpha - E\|_{L^\infty(0, T; X_N(\Omega))} \leq C\alpha, \quad (5)$$

where the constant  $C$  is independent of  $\alpha$  and the set of points  $\{z_j\}_{j=1}^m$  provided that assumption (1) holds.

*Proof.* From (3)-(4), it is obvious that  $E_\alpha - E \in X_N(\Omega)$ , then due to the Green formula we have for any  $\mathbf{v} \in X_N(\Omega)$ :

$$\begin{aligned} \int_{\Omega} \varepsilon_0 \partial_t^2 (E_\alpha - E) \cdot \mathbf{v} \, dx + \int_{\Omega} \frac{1}{\mu_\alpha} \text{curl} (E_\alpha - E) \cdot \text{curl } \mathbf{v} \, dx = \\ \sum_{j=1}^m \left( \frac{1}{\mu_0} - \frac{1}{\mu_j} \right) \int_{z_j + \alpha B_j} \text{curl } E \cdot \text{curl } \mathbf{v} \, dx. \end{aligned} \quad (6)$$

Let  $\mathbf{v}_\alpha$  be defined by

$$\begin{cases} \mathbf{v}_\alpha \in X_N(\Omega), \\ \text{curl } \frac{1}{\mu_\alpha} \text{curl } \mathbf{v}_\alpha = \partial_t(E_\alpha - E) \quad \text{in } \Omega. \end{cases} \quad (7)$$

Then,

$$\begin{aligned} \int_{\Omega} \frac{1}{\mu_\alpha} \text{curl} (E_\alpha - E) \cdot \text{curl } \mathbf{v}_\alpha \, dx &= - \int_{\Omega} \partial_t(E_\alpha - E) \cdot (E_\alpha - E) \, dx = \\ &= - \frac{1}{2} \partial_t \int_{\Omega} |E_\alpha - E|^2 \, dx \end{aligned}$$

and by Green formula, relation (7) gives:

$$\begin{aligned} \int_{\Omega} \partial_t^2 (E_\alpha - E) \cdot \mathbf{v}_\alpha \, dx &= \int_{\Omega} \text{curl } \frac{1}{\mu_\alpha} \text{curl } \partial_t \mathbf{v}_\alpha \cdot \mathbf{v}_\alpha \, dx \\ &= - \int_{\Omega} \frac{1}{\mu_\alpha} \text{curl } \partial_t \mathbf{v}_\alpha \cdot \text{curl } \mathbf{v}_\alpha \, dx \\ &= - \frac{1}{2} \partial_t \int_{\Omega} \frac{1}{\mu_\alpha} |\text{curl } \mathbf{v}_\alpha|^2 \, dx. \end{aligned}$$

Thus, it follows from (6) that

$$\begin{aligned} \varepsilon_0 \partial_t \int_{\Omega} \frac{1}{\mu_{\alpha}} |\operatorname{curl} \mathbf{v}_{\alpha}|^2 dx + \partial_t \int_{\Omega} |E_{\alpha} - E|^2 dx = \\ -2 \sum_{j=1}^m \left( \frac{1}{\mu_0} - \frac{1}{\mu_j} \right) \int_{z_j + \alpha B_j} \operatorname{curl} E \cdot \operatorname{curl} \mathbf{v}_{\alpha} dx. \end{aligned}$$

Next,

$$\left| \sum_{j=1}^m \left( \frac{1}{\mu_0} - \frac{1}{\mu_j} \right) \int_{z_j + \alpha B_j} \operatorname{curl} E \cdot \operatorname{curl} \mathbf{v}_{\alpha} \right| \leq C \| \operatorname{curl} E \|_{\mathbf{L}^2(\mathcal{B}_{\alpha})} \| \operatorname{curl} \mathbf{v}_{\alpha} \|_{\mathbf{L}^2(\Omega)}.$$

Since  $E \in \mathcal{C}^{\infty}([0, T] \times \overline{\Omega})$  we have

$$\| \operatorname{curl} E \|_{\mathbf{L}^2(\mathcal{B}_{\alpha})} \leq \| \operatorname{curl} E \|_{L^{\infty}(\mathcal{B}_{\alpha})} \alpha \left( \sum_{j=1}^m |B_j| \right)^{\frac{1}{2}} \leq C \alpha,$$

which gives

$$\left| \sum_{j=1}^m \left( \frac{1}{\mu_0} - \frac{1}{\mu_j} \right) \int_{z_j + \alpha B_j} \operatorname{curl} E \cdot \operatorname{curl} \mathbf{v}_{\alpha} dx \right| \leq C \alpha \| \operatorname{curl} \mathbf{v}_{\alpha} \|_{\mathbf{L}^2(\Omega)}$$

and so,

$$\varepsilon_0 \partial_t \int_{\Omega} \frac{1}{\mu_{\alpha}} |\operatorname{curl} \mathbf{v}_{\alpha}|^2 dx + \partial_t \int_{\Omega} |E_{\alpha} - E|^2 dx \leq C \alpha \left( \int_{\Omega} \frac{1}{\mu_{\alpha}} |\operatorname{curl} \mathbf{v}_{\alpha}|^2 dx + \int_{\Omega} |E_{\alpha} - E|^2 dx \right)^{1/2}. \quad (8)$$

From the Gronwall Lemma it follows that

$$\left( \int_{\Omega} \frac{1}{\mu_{\alpha}} |\operatorname{curl} \mathbf{v}_{\alpha}|^2 dx \right)^{1/2} + \left( \int_{\Omega} |E_{\alpha} - E|^2 dx \right)^{1/2} \leq C \alpha. \quad (9)$$

Combining this last estimate (9) with the fact that

$$\| \partial_t (E_{\alpha} - E) \|_{L^{\infty}(0, T; H^{-1}(\Omega))} \leq C \| \operatorname{curl} \mathbf{v}_{\alpha} \|_{L^{\infty}(0, T; \mathbf{L}^2(\Omega))}$$

the following estimate holds

$$\| E_{\alpha} - E_0 \|_{L^{\infty}(0, T; \mathbf{L}^2(\Omega))} + \| \partial_t (E_{\alpha} - E_0) \|_{L^{\infty}(0, T; \mathbf{L}^2(\Omega))} \leq C \alpha. \quad (10)$$

Now, taking (formally)  $\mathbf{v} = \partial_t (E_{\alpha} - E)$  in (6) we arrive at

$$\begin{aligned} \varepsilon_0 \partial_t \int_{\Omega} \left[ |\partial_t (E_{\alpha} - E)|^2 + \frac{1}{\mu_{\alpha}} |\operatorname{curl} (E_{\alpha} - E)|^2 \right] dx = \\ 2 \sum_{j=1}^m \left( \frac{1}{\mu_0} - \frac{1}{\mu_j} \right) \int_{z_j + \alpha B_j} \operatorname{curl} E \cdot \operatorname{curl} \partial_t (E_{\alpha} - E) dx. \end{aligned}$$

By using the regularity of  $E$  in  $\Omega$  and estimate (10) given above, we see that

$$\begin{aligned} \left| \sum_{j=1}^m \left( \frac{1}{\mu_0} - \frac{1}{\mu_j} \right) \int_{z_j + \alpha B_j} \operatorname{curl} E \cdot \operatorname{curl} \partial_t (E_\alpha - E) \, dx \right| &\leq C \| \operatorname{curl} E \|_{\mathbf{H}^2(\mathcal{B}_\alpha)} \| \partial_t (E_\alpha - E) \|_{\mathbf{H}^{-1}(\Omega)} \\ &\leq C \alpha^2, \end{aligned}$$

where  $C$  is independent of  $t$  and  $\alpha$ , and so, we obtain

$$\partial_t \int_{\Omega} \left[ |\partial_t (E_\alpha - E)|^2 + \frac{1}{\mu_\alpha} |\operatorname{curl} (E_\alpha - E)|^2 \right] dx \leq C \alpha^2$$

which yields the following estimate

$$\| \partial_t (E_\alpha - E) \|_{L^\infty(0,T;\mathbf{L}^2(\Omega))} + \| E_\alpha - E \|_{L^\infty(0,T;X_N(\Omega))} \leq C \alpha,$$

where  $C$  is independent of  $\alpha$  and the points  $\{z_j\}_{j=1}^m$ .

□

Now, we can estimate  $\operatorname{curl} E_\alpha - \operatorname{curl} E_0$  as follows.

**Proposition 3.1** *Let  $E_\alpha$  and  $E$  be solutions to the problems (3) and (4) respectively. There exist constants  $0 < \alpha_0, C$  such that for  $0 < \alpha < \alpha_0$  the following estimate holds:*

$$\| \operatorname{curl} (E_\alpha - E_0) \|_{L^\infty(0,T;\mathbf{L}^2(\Omega))} \leq C \alpha, \quad (11)$$

*Proof.* To prove estimate (11) it is useful to introduce the following function

$$\hat{v}(x) = \int_0^T v(x, t) z(t) \, dt \in L^2(\Omega), \quad (12)$$

where  $v \in L^1(0, T; L^2(\Omega))$  and  $z(t)$  is a given function in  $\mathcal{C}_0^\infty([0, T])$ .

Then,

$$\hat{E}(x) = \int_0^T E(x, t) z(t) \, dt \text{ and } \hat{E}_\alpha(x) = \int_0^T E_\alpha(x, t) z(t) \, dt \in X(\Omega),$$

which by relation (5) give

$$\begin{cases} (\hat{E}_\alpha - \hat{E}) \in \mathbf{H}^1(\Omega), \\ \operatorname{curl} \operatorname{curl} (\hat{E}_\alpha - \hat{E}) = 0(\alpha) \quad \text{in } \Omega, \\ \operatorname{div} (\hat{E}_\alpha - \hat{E}) = 0 \quad \text{in } \Omega, \\ (\hat{E}_\alpha - \hat{E}) \times \mathbf{n}|_{\partial\Omega} = 0, \end{cases}$$

and so,

$$\| \operatorname{curl} (\hat{E}_\alpha - \hat{E}) \|_{\mathbf{L}^2(\Omega)} = O(\alpha). \quad (13)$$

The fact that  $\operatorname{curl} (E_\alpha - E)$  belongs to  $L^\infty(0, T; \mathbf{L}^2(\Omega))$  and by using estimate (13) we deduce that

$$\int_{\Omega} |\operatorname{curl} E_\alpha(x, t) - \operatorname{curl} E(x, t)|^2 dx = O(\alpha^2) \quad \text{a.e. in } t \in (0, T),$$

which means that

$$\|\operatorname{curl} (E_\alpha - E)\|_{\mathbf{L}^2(\Omega)} = O(\alpha) \quad \text{a.e. in } t \in (0, T).$$

Thus, estimate (11) follows immediately if we take the sup on  $t \in (0, T)$  in the last relation.  $\square$

Before formulating our main result in this section, let us denote  $\Phi_j, j = 1, \dots, m$  the unique vector-valued solution to

$$\begin{cases} \Delta \Phi_j = 0 \text{ in } B_j, \text{ and } \mathbb{R}^d \setminus \overline{B_j}, \\ \Phi_j \text{ is continuous across } \partial B_j, \\ \frac{\mu_j}{\mu_0} \frac{\partial \Phi_j}{\partial \nu_j}|_+ - \frac{\partial \Phi_j}{\partial \nu_j}|_- = -\nu_j, \\ \lim_{|y| \rightarrow +\infty} |\Phi_j(y)| = 0, \end{cases} \quad (14)$$

where  $\nu_j$  denotes the outward unit normal to  $\partial B_j$ , and superscripts  $-$  and  $+$  indicate the limiting values as the point approaches  $\partial B_j$  from outside  $B_j$ , and from inside  $B_j$ , respectively. The existence and uniqueness of this  $\Phi_j$  can be established using single layer potentials with suitably chosen densities, see [6] for the case of conductivity problem. For each inhomogeneity  $z_j + \alpha B_j$  we introduce the polarizability tensor  $M_j$  which is a  $d \times d$ , symmetric, positive definite matrix associated with the  $j$ -th inhomogeneity, given by

$$(M_j)_{k,l} = e_k \cdot \left( \int_{\partial B_j} (\nu_j + \left(\frac{\mu_j}{\mu_0} - 1\right) \frac{\partial \Phi_j}{\partial \nu_j}|_+(y)) y \cdot e_l d\sigma_j(y) \right). \quad (15)$$

Here  $(e_1, \dots, e_d)$  is an orthonormal basis of  $\mathbb{R}^d$ . In terms of this function we are able to prove the following result about the asymptotic behavior of  $\operatorname{curl} E_\alpha \cdot \nu_j|_{\partial(z_j + \alpha B_j)^+}$ .

**Theorem 3.1** *Suppose that (1) is satisfied and let  $\Phi_j, j = 1, \dots, m$  be given as in (14). Then, for the solutions  $E_\alpha, E$  of problems (3) and (4) respectively, and for  $y \in \partial B_j$  we have*

$$\begin{aligned} (\operatorname{curl} E_\alpha(z_j + \alpha y) \cdot \nu_j)|_{\partial(z_j + \alpha B_j)^+} &= \operatorname{curl} E(z_j, t) \cdot \nu_j \\ &+ \left(1 - \frac{\mu_j}{\mu_0}\right) \frac{\partial \Phi_j}{\partial \nu_j}|_+(y) \cdot \operatorname{curl} E(z_j, t) + o(1). \end{aligned} \quad (16)$$

The term  $o(1)$  uniform in  $y \in \partial B_j$  and  $t \in (0, T)$  and depends on the shape of  $\{B_j\}_{j=1}^m$  and  $\Omega$ , the constants  $c_0, T, \mu_0, \{\mu_j\}_{j=1}^m$ , the data  $\varphi, \psi$ , and  $f$ , but is otherwise independent of the points  $\{z_j\}_{j=1}^m$ .

*Proof.*

Let  $\mathcal{E}_\alpha = \operatorname{curl} E_\alpha(x, t)$  and  $\mathcal{E}_0 = \operatorname{curl} E(x, t)$ . Then, according to (3)-(4) we have

$$\varepsilon_0 \partial_t^2 E_\alpha - \operatorname{curl} \frac{1}{\mu_\alpha} \mathcal{E}_\alpha = 0 \text{ and } \operatorname{curl} \mathcal{E}_\alpha = 0, \text{ for } x \in \Omega. \quad (17)$$

We restrict, for simplicity, our attention to the case of a single inhomogeneity, i.e., the case  $m = 1$ . The proof for any fixed number  $m$  of well separated inhomogeneities follows by iteration of the argument that we will present for the case  $m = 1$ . In order to further simplify notation, we assume that the single inhomogeneity has the form  $\alpha B$ , that is, we assume it is centered at the origin. We denote the electromagnetic permeability inside  $\alpha B$  by  $\mu_*$  and define  $\Phi_*$  the same as  $\Phi_j$ , defined in (14), but with  $B_j$  and  $\mu_j$  replaced by  $B$  and  $\mu_*$ , respectively. Define  $\nu$  to be the outward unit normal to  $\partial B$ . Now, following a common practice in multiscale expansions we introduce the local variable  $y = \frac{x}{\alpha}$ , then the domain  $\tilde{\Omega} = (\frac{\Omega}{\alpha})$  is well defined. Next, let  $\varpi$  be given in  $\mathcal{C}_0^\infty([0, T])$ . For any function  $v \in \mathbf{L}^1(0, T; \mathbf{L}^2(\Omega))$ , we define

$$\hat{v}(x) = \int_0^T v(x, t) \varpi(t) dt \in \mathbf{L}^2(\Omega).$$

We remark that  $\widehat{\partial_t v}(x) = - \int_0^T v(x, t) \varpi'(t) dt$ . So that we deduce from (17) that  $\hat{\mathcal{E}}_\alpha$  satisfies

$$\begin{cases} \operatorname{curl} \frac{1}{\mu_\alpha} \hat{\mathcal{E}}_\alpha = \int_0^T E_\alpha \varpi''(t) dt & \text{in } \Omega, \\ \operatorname{curl} \hat{\mathcal{E}}_\alpha = 0 & \text{in } \Omega. \end{cases}$$

Analogously,  $\hat{\mathcal{E}}$  satisfies

$$\begin{cases} \frac{1}{\mu_0} \operatorname{curl} \hat{\mathcal{E}} = \int_0^T E \varpi''(t) dt & \text{in } \Omega, \\ \operatorname{curl} \hat{\mathcal{E}} = 0 & \text{in } \Omega. \end{cases}$$

Indeed, we have  $\hat{\mathcal{E}}_\alpha \times \mathbf{n} = \hat{\mathcal{E}} \times \mathbf{n} = \operatorname{curl}_{\partial\Omega} \hat{f} \times \mathbf{n}$  on the boundary  $\partial\Omega$ , where  $\operatorname{curl}_{\partial\Omega}$  is the tangential curl. Following [4] and [1], we introduce  $q_\alpha^*$  as the unique solution to the following problem

$$\begin{cases} \Delta q_\alpha^* = 0 & \text{in } \tilde{\Omega} = (\frac{\Omega}{\alpha}) \setminus \bar{B} \text{ and in } B, \\ q_\alpha^* \text{ is continuous across } \partial B, \\ \mu_0 \frac{\partial q_\alpha^*}{\partial \nu}|_+ - \mu_* \frac{\partial q_\alpha^*}{\partial \nu}|_- = -(\mu_0 - \mu_*) \hat{\mathcal{E}}(\alpha y) \cdot \nu & \text{on } \partial B, \\ q_\alpha^* = 0 & \text{on } \partial \tilde{\Omega}. \end{cases}$$

The jump condition

$$\mu_0 \frac{\partial q_\alpha^*}{\partial \nu}|_+ - \mu_* \frac{\partial q_\alpha^*}{\partial \nu}|_- = -(\mu_0 - \mu_*) \hat{\mathcal{E}}(\alpha y) \cdot \nu \quad \text{on } \partial B$$

guarantees that  $\hat{\mathcal{E}}_\alpha(x) - \hat{\mathcal{E}}(x) - \text{grad}_y q_\alpha^*(\frac{x}{\alpha})$  belongs to the functional space  $X_N(\Omega)$ , where  $\text{grad}_{\partial\Omega}$  is the tangential gradient. Since

$$\begin{cases} \text{curl} \frac{1}{\mu_\alpha} (\hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{grad}_y q_\alpha^*(\frac{x}{\alpha})) = \int_0^T \left[ E_\alpha - \chi(\Omega \setminus \overline{\alpha B}) E + \frac{\mu_*}{\mu_0} \chi(\alpha B) E \right] \varpi''(t) dt & \text{in } \Omega, \\ \text{curl} (\hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{grad}_y q_\alpha^*(\frac{x}{\alpha})) = 0 & \text{in } \Omega, \\ (\hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{grad}_y q_\alpha^*(\frac{x}{\alpha})) \times \mathbf{n} = 0 & \text{on } \partial\Omega, \end{cases}$$

where  $\chi(\omega)$  is the characteristic function of the domain  $\omega$ , we arrive, as a consequence of the energy estimate given by Lemma 3.1, at the following

$$\begin{cases} (\hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{grad}_y q_\alpha^*(\frac{x}{\alpha})) \in X_N(\Omega), \\ \text{curl} \frac{1}{\mu_\alpha} (\hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{curl}_y q_\alpha^*(\frac{x}{\alpha})) = 0(\alpha) & \text{in } \Omega, \\ \text{curl} (\hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{grad}_y q_\alpha^*(\frac{x}{\alpha})) = 0 & \text{in } \Omega, \\ (\hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{grad}_y q_\alpha^*(\frac{x}{\alpha})) \times \mathbf{n} = 0 & \text{on } \partial\Omega. \end{cases}$$

From [4] we know that this yields the following estimate

$$\| \text{curl} \frac{1}{\mu_\alpha} (\hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{grad}_y q_\alpha^*(\frac{x}{\alpha})) \|_{L^2(\Omega)} + \| \hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{grad}_y q_\alpha^*(\frac{x}{\alpha}) \|_{L^2(\Omega)} \leq C\alpha,$$

and so,

$$(\hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{grad}_y q_\alpha^*(\frac{x}{\alpha})) \cdot \nu|_+ = 0(\alpha) \quad \text{on } \partial(\alpha B).$$

Now, we denote by  $q_*$  be the unique (scalar) solution to

$$\begin{cases} \Delta q_* = 0 & \text{in } \mathbb{R}^d \setminus \overline{B} \text{ and in } B, \\ q_* \text{ is continuous across } \partial B, \\ \mu_0 \frac{\partial q_*}{\partial \nu}|_+ - \mu_* \frac{\partial q_*}{\partial \nu}|_- = -(\mu_0 - \mu_*) \hat{\mathcal{E}}(0) \cdot \nu & \text{on } \partial B, \\ \lim_{|y| \rightarrow +\infty} q_* = 0. \end{cases}$$

In the spirit of Theorem 1 in [6] it follows that

$$\| (\text{grad}_y q_* - \text{grad}_y q_\alpha^*)(\frac{x}{\alpha}) \|_{L^2(\Omega)} \leq C\alpha^{1/2},$$

which yields

$$(\hat{\mathcal{E}}_\alpha - \hat{\mathcal{E}} - \text{grad}_y q_*(\frac{x}{\alpha})) \cdot \nu = o(1) \quad \text{on } \partial(\alpha B).$$

Writing  $q_*$  in terms of  $\Phi_*$  gives

$$\int_0^T \left[ (\text{curl } E_\alpha(\alpha y) \cdot \nu)|_{\partial(\alpha B)^+} - \nu \cdot \text{curl } E(0, t) - \left( \frac{\mu_0}{\mu_*} - 1 \right) \frac{\partial \Phi_*}{\partial \nu}|_+(y) \cdot \text{curl } E(0, t) \right] \varpi(t) dt = o(1),$$

for any  $\varpi \in \mathcal{C}_0^\infty([0, T])$ , and so, by iterating the same argument for the case of  $m$  (well separated) inhomogeneities  $z_j + \alpha B_j, j = 1, \dots, m$ , we arrive at the promised asymptotic formula (16).  $\square$

#### 4 The identification procedure

Before describing our identification procedure, let us introduce the following cutoff function  $\beta(x) \in \mathcal{C}_0^\infty(\Omega)$  such that  $\beta \equiv 1$  in a subdomain  $\Omega'$  of  $\Omega$  that contains the inhomogeneities  $\mathcal{B}_\alpha$  and let  $\eta \in \mathbb{R}^d$ . We will take in what follows  $E(x, t) = \eta^\perp e^{i\eta \cdot x - i\sqrt{\mu_0}|\eta|t}$  where  $\eta^\perp$  is a unit vector that is orthogonal to  $\eta$  which corresponds to taking  $\varphi(x) = \eta^\perp e^{i\eta \cdot x}$ ,  $\psi(x) = -i\sqrt{\mu_0}|\eta|\eta^\perp e^{i\eta \cdot x}$ , and  $f(x, t) = \eta^\perp \times \mathbf{n} e^{i\eta \cdot x - i\sqrt{\mu_0}|\eta|t}$  and assume that we are in possession of the measurements of:

$$\operatorname{curl} E_\alpha \times \mathbf{n} \quad \text{on } \Gamma \times (0, T),$$

where  $\Gamma$  is an open part of  $\partial\Omega$ . Suppose now that  $T$  and the part  $\Gamma$  of the boundary  $\partial\Omega$  are such that they geometrically control  $\Omega$  which roughly means that every geometrical optic ray, starting at any point  $x \in \Omega$  at time  $t = 0$  hits  $\Gamma$  before time  $T$  at a non diffractive point, see [5]. It follows from [17] (see also [13], [11] and [12]) that there exists (a unique)  $g_\eta \in H_0^1(0, T; TL^2(\Gamma))$  (constructed by the Hilbert Uniqueness Method) such that the unique weak solution  $w_\eta$  to

$$\begin{cases} (\partial_t^2 + \operatorname{curl} \operatorname{curl}) w_\eta = 0 & \text{in } \Omega \times (0, T), \\ \operatorname{div} w_\eta = 0 & \text{in } \Omega \times (0, T), \\ w_\eta|_{t=0} = \beta(x)\eta^\perp e^{i\eta \cdot x}, \partial_t w_\eta|_{t=0} = 0 & \text{in } \Omega, \\ w_\eta \times \mathbf{n}|_{\partial\Omega \setminus \Gamma \times (0, T)} = 0, \\ w_\eta \times \mathbf{n}|_{\Gamma \times (0, T)} = g_\eta, \end{cases} \quad (18)$$

satisfies  $w_\eta(T) = \partial_t w_\eta(T) = 0$  in  $\Omega$ .

Let  $\theta_\eta \in H^1(0, T; TL^2(\Gamma))$  denote the unique solution of the Volterra equation of second kind

$$\begin{cases} \partial_t \theta_\eta(x, t) + \int_t^T e^{-i|\eta|(s-t)} (\theta_\eta(x, s) - i|\eta| \partial_t \theta_\eta(x, s)) ds = g_\eta(x, t) & \text{for } x \in \Gamma, t \in (0, T), \\ \theta_\eta(x, 0) = 0 & \text{for } x \in \Gamma. \end{cases} \quad (19)$$

The existence and uniqueness of this  $\theta_\eta$  in  $H^1(0, T; TL^2(\Gamma))$  for any  $\eta \in \mathbb{R}^d$  can be established using the resolvent kernel. However, observing from differentiation of (19) with respect to  $t$  that  $\theta_\eta$  is the unique solution of the ODE:

$$\begin{cases} \partial_t^2 \theta_\eta - \theta_\eta = e^{i|\eta|t} \partial_t (e^{-i|\eta|t} g_\eta) & \text{for } x \in \Gamma, t \in (0, T), \\ \theta_\eta(x, 0) = 0, \partial_t \theta_\eta(x, T) = 0 & \text{for } x \in \Gamma, \end{cases} \quad (20)$$

the function  $\theta_\eta$  may be find (in practice) explicitly with variation of parameters and it also immediately follows from this observation that  $\theta_\eta$  belongs to  $H^2(0, T; TL^2(\Gamma))$ .

We introduce  $v_\eta$  as the unique weak solution (obtained by transposition as done in [15] and in [14] [Theorem 4.2, page 46] for the scalar function) in  $\mathcal{C}^0(0, T; X(\Omega)) \cap \mathcal{C}^1(0, T; L^2(\Omega))$  to the



following problem

$$\left\{ \begin{array}{l} (\partial_t^2 + \operatorname{curl} \operatorname{curl}) v_\eta = 0 \quad \text{in } \Omega \times (0, T), \\ \operatorname{div} v_\eta = 0 \quad \text{in } \Omega \times (0, T), \\ v_\eta|_{t=0} = 0 \quad \text{in } \Omega, \\ \partial_t v_\eta|_{t=0} = \sum_{j=1}^m i(1 - \frac{\mu_0}{\mu_j}) \eta \times (\nu_j + (\frac{\mu_0}{\mu_j} - 1) \frac{\partial \Phi_j}{\partial \nu_j}|_+) e^{i\eta \cdot z_j} \delta_{\partial(z_j + \alpha B_j)} \in Y(\Omega) \quad \text{in } \Omega, \\ v_\eta \times \mathbf{n}|_{\partial\Omega \times (0, T)} = 0. \end{array} \right.$$

Then, the following holds.

**Proposition 4.1** Suppose that  $\Gamma$  and  $T$  geometrically control  $\Omega$ . For any  $\eta \in \mathbb{R}^d$  we have

$$\begin{aligned} \int_0^T \int_\Gamma g_\eta \cdot (\operatorname{curl} v_\eta \times \mathbf{n}) \, d\sigma(x) dt &= \alpha^2 \sum_{j=1}^m \mu_0 (1 - \frac{\mu_j}{\mu_0}) e^{2i\eta \cdot z_j} \eta \cdot \int_{\partial B_j} (\nu_j \\ &+ (\frac{\mu_j}{\mu_0} - 1) \frac{\partial \Phi_j}{\partial \nu_j}|_+(y)) \eta \cdot y \, ds_j(y) + o(\alpha^2). \end{aligned} \quad (21)$$

*Proof.* Multiply the equation  $(\partial_t^2 + \operatorname{curl} \operatorname{curl}) v_\eta = 0$  by  $w_\eta$  and integrating by parts over  $(0, T) \times \Omega$ , for any  $\eta \in \mathbb{R}^d$  we have

$$\begin{aligned} \alpha \sum_{j=1}^m i(1 - \frac{\mu_j}{\mu_0}) e^{2i\eta \cdot z_j} \eta \cdot \int_{\partial B_j} (\nu_j + (\frac{\mu_j}{\mu_0} - 1) \frac{\partial \Phi_j}{\partial \nu_j}|_+(y)) e^{i\alpha \eta \cdot y} \, ds(y) = \\ -\mu_0^{-1} \int_0^T \int_\Gamma g_\eta \cdot (\operatorname{curl} v_\eta \times \mathbf{n}) \, d\sigma(x) dt. \end{aligned}$$

Now, we take the Taylor expansion of  $\alpha e^{i\alpha \eta \cdot y}$  in the left side of the last equation, we obtain the convenient asymptotic formula (21).  $\square$

To identify the locations and certain properties of the small inhomogeneities  $\mathcal{B}_\alpha$  let us view the averaging of the boundary measurements

$$\operatorname{curl} E_\alpha \times \mathbf{n}|_{\Gamma \times (0, T)},$$

using the solution  $\theta_\eta$  to the Volterra equation (19) or equivalently the ODE (20), as a function of  $\eta$ . The following holds.

**Theorem 4.1** Let  $\eta \in \mathbb{R}^d$ . Let  $E_\alpha$  be the unique solution in  $\mathcal{C}^0(0, T; X(\Omega)) \cap \mathcal{C}^1(0, T; L^2(\Omega))$  to the Maxwell's equations (3) with  $\varphi(x) = \eta^\perp e^{i\eta \cdot x}$ ,  $\psi(x) = -i\sqrt{\mu_0}|\eta|\eta^\perp e^{i\eta \cdot x}$ , and  $f(x, t) = \eta^\perp e^{i\eta \cdot x - i\sqrt{\mu_0}|\eta|t}$ . Suppose that  $\Gamma$  and  $T$  geometrically control  $\Omega$ , then we have

$$\begin{aligned} & \int_0^T \int_\Gamma \left[ \theta_\eta \cdot (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}) + \partial_t \theta_\eta \cdot \partial_t (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}) \right] d\sigma(x) dt = \\ & \alpha^2 \sum_{j=1}^m (\mu_0 - \mu_j) e^{2i\eta \cdot z_j} M_j(\eta) \cdot \eta + O(\alpha^2), \end{aligned} \quad (22)$$

where  $\theta_\eta$  is the unique solution to the Volterra equation (20) with  $g_\eta$  defined as the boundary control in (18) and  $M_j$  is the polarization tensor of  $B_j$ , defined by

$$(M_j)_{k,l} = e_k \cdot \left( \int_{\partial B_j} \left( \nu_j + \left( \frac{\mu_j}{\mu_0} - 1 \right) \frac{\partial \Phi_j}{\partial \nu_j} \Big|_+(y) \right) y \cdot e_l ds_j(y) \right). \quad (23)$$

Here  $(e_1, e_2)$  is an orthonormal basis of  $\mathbb{R}^d$ . The term  $O(\alpha^2)$  is independent of the points  $\{z_j, j = 1, \dots, m\}$ .

*Proof.* From  $\partial_t \theta_\eta(T) = 0$  and  $(\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n})|_{t=0} = 0$  the term  $\int_0^T \int_\Gamma \partial_t \theta_\eta \cdot \partial_t (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}) d\sigma(x) dt$  has to be interpreted as follows

$$\int_0^T \int_\Gamma \partial_t \theta_\eta \cdot \partial_t (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}) = - \int_0^T \int_\Gamma \partial_t^2 \theta_\eta \cdot (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}). \quad (24)$$

Next, introduce

$$\tilde{E}_{\alpha,\eta}(x, t) = E(x, t) + \int_0^t e^{-i\sqrt{\mu_0}|\eta|s} v_\eta(x, t-s) ds, x \in \Omega, t \in (0, T). \quad (25)$$

We have

$$\begin{aligned} & \int_0^T \int_\Gamma \left[ \theta_\eta \cdot (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}) + \partial_t \theta_\eta \cdot \partial_t (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}) \right] = \\ & \int_0^T \int_\Gamma \left[ \theta_\eta \cdot (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} \tilde{E}_{\alpha,\eta} \times \mathbf{n}) + \partial_t \theta_\eta \cdot \partial_t (\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} \tilde{E}_{\alpha,\eta} \times \mathbf{n}) \right] \\ & + \int_0^T \int_\Gamma \left[ \theta_\eta \cdot \int_0^t e^{-i\sqrt{\mu_0}|\eta|s} v_\eta(x, t-s) \times \mathbf{n} ds + \partial_t \theta_\eta \cdot \partial_t \int_0^t e^{-i\sqrt{\mu_0}|\eta|s} v_\eta(x, t-s) \times \mathbf{n} ds \right]. \end{aligned}$$

Since  $\theta_\eta$  satisfies the Volterra equation (20) and

$$\begin{aligned} & \partial_t \left( \int_0^t e^{-i\sqrt{\mu_0}|\eta|s} v_\eta(x, t-s) \times \mathbf{n} ds \right) = \partial_t \left( -e^{-i\sqrt{\mu_0}|\eta|t} \int_0^t e^{i\sqrt{\mu_0}|\eta|s} v_\eta(x, s) \times \mathbf{n} ds \right) \\ & = i\sqrt{\mu_0}|\eta| e^{-i\sqrt{\mu_0}|\eta|t} \int_0^t e^{i\sqrt{\mu_0}|\eta|s} v_\eta(x, s) \times \mathbf{n} ds + v_\eta(x, t) \times \mathbf{n}, \end{aligned}$$

we obtain by integrating by parts over  $(0, T)$  that

$$\begin{aligned} & \int_0^T \int_{\Gamma} \left[ \theta_{\eta} \cdot \int_0^t e^{-i\sqrt{\mu_0}|\eta|s} v_{\eta}(x, t-s) \times \mathbf{n} \, ds + \partial_t \theta_{\eta} \cdot \partial_t \int_0^t e^{-i\sqrt{\mu_0}|\eta|s} v_{\eta}(x, t-s) \times \mathbf{n} \, ds \right] \\ &= \int_0^T \int_{\Gamma} (v_{\eta}(x, t) \times \mathbf{n}) \cdot (\partial_t \theta_{\eta} + \int_t^T \theta_{\eta}(s) e^{i\sqrt{\mu_0}|\eta|(t-s)} \, ds) \\ & \quad - i\sqrt{\mu_0}|\eta| (e^{-i\sqrt{\mu_0}|\eta|t} \partial_t \theta_{\eta}(t)) \cdot \int_0^t e^{i\sqrt{\mu_0}|\eta|s} v_{\eta}(x, s) \times \mathbf{n} \, ds \, dt \\ &= \int_0^T \int_{\Gamma} v_{\eta}(x, t) \times \mathbf{n} \cdot (\partial_t \theta_{\eta} + \int_t^T (\theta_{\eta}(s) - i\sqrt{\mu_0}|\eta| \partial_t \theta_{\eta}(s)) e^{i\sqrt{\mu_0}|\eta|(t-s)} \, ds) \, dt \\ &= \int_0^T \int_{\Gamma} g_{\eta}(x, t) \cdot (\operatorname{curl} v_{\eta}(x, t) \times \mathbf{n}) \, dt \end{aligned}$$

and so, from Proposition 4.1 we obtain

$$\begin{aligned} & \int_0^T \int_{\Gamma} \left[ \theta_{\eta} \cdot (\operatorname{curl} E_{\alpha} \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}) + \partial_t \theta_{\eta} \cdot \partial_t (\operatorname{curl} E_{\alpha} \times \mathbf{n} - \operatorname{curl} E \times \mathbf{n}) \right] = \\ & \alpha^2 \sum_{j=1}^m \left( 1 - \frac{\mu_j}{\mu_0} \right) e^{2i\eta \cdot z_j} \eta \cdot \int_{\partial B_j} \left( \nu_j + \left( \frac{\mu_j}{\mu_0} - 1 \right) \frac{\partial \Phi_j}{\partial \nu_j} \Big|_+(y) \right) \eta \cdot y \, ds_j(y) \\ & + \int_0^T \int_{\Gamma} \left[ \theta_{\eta} \cdot (\operatorname{curl} E_{\alpha} \times \mathbf{n} - \operatorname{curl} \tilde{E}_{\alpha, \eta} \times \mathbf{n}) + \partial_t \theta_{\eta} \cdot \partial_t (\operatorname{curl} E_{\alpha} \times \mathbf{n} \right. \\ & \quad \left. - \operatorname{curl} \tilde{E}_{\alpha, \eta} \times \mathbf{n}) \right] + o(\alpha^2). \end{aligned}$$

In order to prove Theorem 4.1 it suffices then to show that

$$\int_0^T \int_{\Gamma} \left[ \theta_{\eta} \cdot (\operatorname{curl} E_{\alpha} \times \mathbf{n} - \operatorname{curl} \tilde{E}_{\alpha, \eta} \times \mathbf{n}) + \partial_t \theta_{\eta} \cdot \partial_t (\operatorname{curl} E_{\alpha} \times \mathbf{n} - \operatorname{curl} \tilde{E}_{\alpha, \eta} \times \mathbf{n}) \right] = o(\alpha^2). \quad (26)$$

Since

$$\left\{ \begin{aligned} & (\partial_t^2 - \operatorname{curl} \frac{1}{\mu_0} \operatorname{curl}) \left( \int_0^t e^{-i\sqrt{\mu_0}|\eta|s} v_{\eta}(x, t-s) \, ds \right) \\ &= \sum_{j=1}^m i \left( 1 - \frac{\mu_j}{\mu_0} \right) \eta \times \left( \nu_j + \left( \frac{\mu_j}{\mu_0} - 1 \right) \frac{\partial \Phi_j}{\partial \nu_j} \Big|_+(y) \right) e^{i\eta \cdot z_j} \delta_{\partial(z_j + \alpha B_j)} e^{-i\sqrt{\mu_0}|\eta|t} \quad \text{in } \Omega \times (0, T), \\ & \left( \int_0^t e^{-i\sqrt{\mu_0}|\eta|s} v_{\eta}(x, t-s) \, ds \right) \Big|_{t=0} = 0, \partial_t \left( \int_0^t e^{-i\sqrt{\mu_0}|\eta|s} v_{\eta}(x, t-s) \, ds \right) \Big|_{t=0} = 0 \quad \text{in } \Omega, \\ & \left( \int_0^t e^{-i\sqrt{\mu_0}|\eta|s} v_{\eta}(x, t-s) \, ds \right) \times \mathbf{n} \Big|_{\partial \Omega \times (0, T)} = 0, \end{aligned} \right.$$

it follows from Theorem 3.1 that

$$\left\{ \begin{aligned} & (\partial_t^2 - \operatorname{curl} \frac{1}{\mu_0} \operatorname{curl}) (E_{\alpha} - \tilde{E}_{\alpha, \eta}) = o(\alpha^2) \quad \text{in } \Omega \times (0, T), \\ & (E_{\alpha} - \tilde{E}_{\alpha, \eta}) \Big|_{t=0} = 0, \partial_t (E_{\alpha} - \tilde{E}_{\alpha, \eta}) \Big|_{t=0} = 0 \quad \text{in } \Omega, \\ & (E_{\alpha} - \tilde{E}_{\alpha, \eta}) \times \mathbf{n} \Big|_{\partial \Omega \times (0, T)} = 0. \end{aligned} \right.$$

Following the proof of Proposition 3.1, we immediately obtain

$$\|E_\alpha - \tilde{E}_{\alpha,\eta}\|_{L^2(\Omega)} = o(\alpha^2), t \in (0, T), x \in \Omega,$$

where  $o(\alpha^2)$  is independent of the points  $\{z_j\}_{j=1}^m$ . To prove (26) it suffices then from (24) to show that the following estimate holds

$$\|\operatorname{curl} E_\alpha \times \mathbf{n} - \operatorname{curl} \tilde{E}_{\alpha,\eta} \times \mathbf{n}\|_{L^2(0,T;TL^2(\Gamma))} = o(\alpha^2).$$

Let  $\theta$  be given in  $\mathcal{C}_0^\infty(]0, T[)$  and define

$$\hat{\tilde{E}}_{\alpha,\eta}(x) = \int_0^T \tilde{E}_{\alpha,\eta}(x, t)\theta(t) dt$$

and

$$\hat{E}_\alpha(x) = \int_0^T E_\alpha(x, t)\theta(t) dt.$$

From definition (25) we can write

$$\left\{ \begin{array}{l} (\hat{E}_\alpha - \hat{\tilde{E}}_\alpha) \in \mathbf{H}^1(\Omega), \\ \operatorname{curl} \operatorname{curl} (\hat{E}_\alpha - \hat{\tilde{E}}_\alpha) = 0(\alpha) \in Y(\Omega) \quad \text{in } \Omega, \\ \operatorname{div} (\hat{E}_\alpha - \hat{\tilde{E}}_\alpha) = 0 \quad \text{in } \Omega, \\ (\hat{E}_\alpha - \hat{\tilde{E}}_\alpha) \times \mathbf{n}|_{\partial\Omega} = 0. \end{array} \right. \quad (27)$$

In the spirit of the standard elliptic regularity [9] we deduce for the boundary value problem (27) that

$$\|\operatorname{curl} (\hat{E}_\alpha - \hat{\tilde{E}}_\alpha) \times \mathbf{n}\|_{\mathbf{L}^2(\Gamma)} = O(\alpha^2),$$

for all  $\theta \in \mathcal{C}_0^\infty(]0, T[)$ ; whence

$$\|\operatorname{curl} (E_\alpha - \hat{E}_\alpha) \times \mathbf{n}\|_{\mathbf{L}^2(\Gamma)} = o(\alpha^2) \text{ a. e. in } t \in (0, T),$$

and so, the desired estimate (22) holds. The proof of Theorem 4.1 is then over.  $\square$

Our identification procedure is deeply based on Theorem 4.1. Let us neglect the asymptotically small remainder in the asymptotic formula (22), and define  $\aleph_\alpha(\eta)$  by

$$\aleph_\alpha(\eta) = \int_0^T \int_\Gamma \left[ \theta_\eta \cdot (\operatorname{curl} (E_\alpha - E) \times \mathbf{n}) + \partial_t \theta_\eta \cdot \partial_t (\operatorname{curl} (E_\alpha - E) \times \mathbf{n}) \right].$$

Recall that the function  $e^{2i\eta \cdot z_j}$  is exactly the Fourier Transform (up to a multiplicative constant) of the Dirac function  $\delta_{-2z_j}$  (a point mass located at  $-2z_j$ ). From Theorem 4.1 it follows that the function  $e^{2i\eta \cdot z_j}$  is (approximately) the Fourier Transform of a linear combination of derivatives of point masses, or

$$\check{\aleph}_\alpha(\eta) \approx \alpha^2 \sum_{j=1}^m L_j \delta_{-2z_j},$$

where  $L_j$  is a second order constant coefficient, differential operator whose coefficients depend on the polarization tensor  $M_j$  defined by (23) (see [6] for its properties) and  $\check{\aleph}_\alpha(\eta)$  represents the inverse Fourier Transform of  $\aleph_\alpha(\eta)$ . The reader is referred to [6] for properties of the tensor polarization  $M_j$ .

The method of reconstruction consists in sampling values of  $\check{\aleph}_\alpha(\eta)$  at some discrete set of points and then calculating the corresponding discrete inverse Fourier Transform. After a rescaling the support of this discrete inverse Fourier Transform yields the location of the small inhomogeneities  $\mathcal{B}_\alpha$ . Once the locations are known we may calculate the polarization tensors  $(M_j)_{j=1}^m$  by solving an appropriate linear system arising from (22). This procedure generalizes the approach developed in [3] for the two-dimensional (time-independent) inverse conductivity problem and generalize the results in [1] to the full time-dependent Maxwell's equations.

## 5 Conclusion

In this paper, we are convinced that the use of approximate formulae such as (22) represents a very promising approach to the dynamical identification of small inhomogeneities that are embedded in a homogeneous medium. We also believe that our method yields a good approximation to small amplitude perturbations in the electromagnetic parameters (for the example of electric permittivity  $\varepsilon_\alpha(x) = \varepsilon_0 + \alpha\varepsilon(x)$ ) from the measurements:

$$\operatorname{curl} H_\alpha \times \mathbf{n} \quad \text{on } \Gamma \times (0, T).$$

Our method may yield the Fourier transform of the amplitude perturbation  $\varepsilon(x)$ . This issue will be considered in a forthcoming work [7].

## References

- [1] H. AMMARI, *An inverse initial boundary value problem for the wave equation in the presence of imperfections of small volume*. SIAM J. Control Optim. 41 (2003), 1194-1211.
- [2] H. AMMARI, H. KANG, E. KIM, and M. LIM, *Reconstruction of closely spaced small inclusions*. SIAM J. Numer. Anal. 42, no. 6, (2005), 2408-2428.
- [3] H. AMMARI, S. MOSKOW, and M. VOGELIUS, *Boundary integral formulas for the reconstruction of electromagnetic imperfections of small diameter*, ESAIM: COCV 9 (2003), 49-66.
- [4] H. AMMARI, M. VOGELIUS, and D. VOLKOV, *Asymptotic formulas for perturbations in the electromagnetic fields due to the presence of inhomogeneities of small diameter II. The full Maxwell equations*, J. Math. Pures Appl. 80 (2001), 769-814.
- [5] C. BARDOS, G. LEBEAU, and J. RAUCH, *Sharp sufficient conditions for the observation, control and stabilization of waves from the boundary*, SIAM J. Control Opt. 30 (1992), 1024-1065.
- [6] D. J. CEDIO-FENGYA, S. MOSKOW, and M. VOGELIUS, *Identification of conductivity imperfections of small diameter by boundary measurements. Continuous dependence and computational reconstruction*, Inverse Problems 14 (1998), 553-595.

- [7] C. DAVEAU, and A. KHELIFI, *An inverse problem for the time dependent Maxwell's equations in the presence of imperfections of small volume*, preprint.
- [8] C. DAVEAU, A. ZAGHDANI, *Two Inequalities of Poincare-Friedrichs on discontinuous spaces for Maxwell's equations*, C. R. Acad. Sci. Paris, Ser. I 342, (2006).
- [9] L. C. EVANS, *Partial Differential Equations*, Graduate Studies in Mathematics, AMS, 1998, Providence, Rhode Island.
- [10] A. FRIEDMAN and M. VOGELIUS, *Identification of small inhomogeneities of extreme conductivity by boundary measurements: a theorem on continuous dependence*, Arch. Rat. Mech. Anal. 105 (1989), 299-326.
- [11] K. A. KIME, *Boundary controllability of Maxwell's equations in a spherical region*, SIAM J. Control Optim. 28 (1990), 294-319.
- [12] V. KOMORNIK, *Boundary stabilization, observation and control of Maxwell's equations*, Panamer. Math. J. 4 (1994), 47-61.
- [13] J. E. LAGNESE, *Exact boundary controllability of Maxwell's equations in a general region*, SIAM J. Control Optim. 27 (1989), 374-388.
- [14] J. L. LIONS, *Contrôlabilité exacte, Perturbations et Stabilisation de Systèmes Distribués, Tome 1, Contrôlabilité Exacte*, Masson 1988, Paris.
- [15] J. L. LIONS and E. MAGENES, *Nonhomogeneous Boundary Value Problems and Applications*, Vol. 1, Springer, 1972.
- [16] S. R. McDOWALL, *An electromagnetic inverse problem in chiral media*, Trans. Amer. Math. Soc. 352 (2000), 2993-3013.
- [17] S. NICAISE, *Exact boundary controllability of Maxwell's equations in heterogeneous media and an application to an inverse source problem*, SIAM J. Control Optim. 38 (2000), 1145-1170.
- [18] V. G. ROMANOV and S. I. KABANIKHIN, *Inverse Problems for Maxwell's Equations*, Inverse and Ill-posed Problems Series, VSP, Utrecht, 1994.
- [19] J. SYLVESTER and G. UHLMANN, *A global uniqueness theorem for an inverse boundary value problem*, Ann. Math. 125 (1987), 153-169.
- [20] E. SOMERSALO, D. ISAACSON, and M. CHENEY, *A linearized inverse boundary value problem for Maxwell's equations*, J. Comput. Appl. Math. 42 (1992), 123-136.
- [21] M. VOGELIUS and D. VOLKOV, *Asymptotic formulas for perturbations in the electromagnetic fields due to the presence of inhomogeneities*, Math. Model. Numer. Anal. 34 (2000), 723-748.
- [22] D. VOLKOV, *An inverse problem for the time harmonic Maxwell equations*, PhD thesis, Rutgers University, New Brunswick, NJ, 2001.

## **Current address**

### **Christian Daveau**

Département de Mathématiques, Site Saint-Martin II, BP 222, & Université de Cergy-Pontoise, 95302 Cergy-Pontoise Cedex, France e-mail: christian.daveau@math.u-cergy.fr).

**KHELIFI Abdessatar**

Département de Mathématiques & Informatique Faculté des Sciences, 7021 Zarzouna - Bizerte, Tunisia e-mail:abdessatar.khelifi@fsb.rnu.tn

**SUSHCHENKO Anton**

ETIS & UMR CNRS 8051, 6 avenue du Ponceau, BP 44, 95014 Cergy-Pontoise Cedex, France (Email: anton.sushchenko@ensea.fr).





## SHARK-FISH INTERPLAY AT DIFFERENT LIFESTAGES

CHATTERJEE Samrat, (IND), VENTURINO Ezio, (I)

**Abstract.** An LSF (larvae-shark-fish) food chain is described and analyzed, under reasonable natural assumptions, to assess the possible control the intermediate element in the chain can exert on their top predators. Some analytical results are reported. In view of the complexity of the coexistence equilibrium, for understanding the system's behavior in its surroundings, numerical simulations are performed. A particular case is then examined in more detail. Stability of the coexistence equilibrium as well as the one of the larvae and shark-free equilibrium is studied. An interpretation of the results in ecological terms is provided. We conclude the paper with some considerations on the role of fisheries in this aquatic ecosystem.

**Key words and phrases.** predator-prey models, food webs, equilibria, stability, alternative food source.

*Mathematics Subject Classification.* Primary 92D25, 92D40; Secondary 37N25.

### 1 Introduction

Mathematical population theory studies interactions of different populations which share at least partially a common environment. It originates from the researches of Volterra, D'Ancona and Lotka in the early twentieth century, who combined the models of Malthus and Verhulst of the previous century to describe interactions among predators and prey in nature. The classical Lotka-Volterra model for fish interactions accounts for the oscillations found in the data of fisheries in the Adriatic sea in the years following World War I. The model, due to an intrinsic bad feature, namely the neutral stability of its equilibrium point, has been later revised and improved. Studies in mathematical biology have then continued along the century and have witnessed an upsurge in the seventies.

More complex models have been investigated, most notably food webs, in which in general a top predator feeds on some prey, which themselves feed on other organisms at a lower trophic

level. This can continue through several levels, until a bottom prey is found. In the ocean, the latter is found at the planktonic level, in particular it is given by phytoplankton, on which zooplankton itself feeds.

In this paper we consider a “low” level food chain, made only of three levels, but in which the bottom one is made by the younglings of the top predators. In an aquatic environment, the latter are at the larval stage. In the environment there are also fishes present, as well as the adult predators. The larvae would mature into adults, if they survive the attacks of other fishes. In fact they constitute a food source for all the other fishes, which in turn are predated upon by the adult sharks. In this way, the interaction of fish with the sharks larvae might represent an indirect control on the sharks themselves. With this model we focus on the relationships among these populations and highlight some of their consequences.

Classical researches in field show that food webs are structured via complex interactions between consumers, by the top down approach, while the from the bottom up approach shows the role of resources, [2, 3, 6, 7, 8, 9].

## 2 The mathematical model description

We consider the ocean environment, in which the large predators, like the sharks, feed on all other fishes. The latter have possibly different food sources, like for instance plankton, but can also feed on the larvae of the former, thereby possibly reducing the number of top predators. Let therefore  $L$  denote the sharks at larval stage,  $S$  denote the adult sharks, able to reproduce, and  $F$  denote the fishes which act as predators on the larvae and as prey for the sharks. Note also that the latter do not feed on their own larvae, as these are basically disregarded in view of their too small size. The resulting model reads as follows

$$\begin{aligned}\frac{dL}{dt} &= -nL + a\tilde{c}SF - bLF \\ \frac{dS}{dt} &= gL - mS \\ \frac{dF}{dt} &= rF \left(1 - \frac{F}{K}\right) + bcLF - aSF.\end{aligned}\tag{1}$$

Here, several parameters come into play. We denote by  $n$  the loss in the larval population due both to the maturation process and natural deaths,  $g \leq n$  represents the maturity rate of the sharks from larval to adult,  $m$  is the natural death rates of the adult sharks. By  $r$  we indicate the specific net growth rate of the fishes, which may be positive if their birth rate is greater than their death rate and food sources other than larvae are available. It may also be negative in case the fishes’ death rate exceeds the birth rate. The rate at which fishes predate on the sharks larvae is  $b$  while  $a$  is the corresponding hunting rate of adult sharks predated on the fishes. For  $r > 0$ ,  $K$  represents the fish environment carrying capacity. For mathematical simplicity, we assume here that the conversion rates  $c$  and  $\tilde{c}$  due to the predation is equal to 1 for both fish and sharks.

The first equation describes the evolution of the larvae, which are born from the adult sharks when they can feed on the fish, second term, and are removed by predation, third term, or once

they die or become adult, first term. The second equation gives the sharks dynamics: they enter this class via larvae maturation, first term and are subject only to natural mortality not having higher predators in the food chain. The fishes, third equation, experience possibly logistic growth, first two terms, also predate on larvae, second term, and are hunted by predators, third term.

## 2.1 A boundedness result for the trajectories of the dynamical system

We consider  $Z$ , the total environment population, a function of time,

$$Z = L + S + F \quad (2)$$

and take the time-derivative of (2) along the solutions of (1), thus obtaining

$$\frac{dZ}{dt} = (g - n)L - mS + rF \left(1 - \frac{F}{K}\right) \leq -nL - mS + rF \left(1 - \frac{F}{K}\right).$$

Taking an arbitrary constant  $\eta > 0$  we get,

$$\frac{dZ}{dt} + \eta Z \leq \eta Z - nL - mS + rF \left(1 - \frac{F}{K}\right).$$

Now if we choose  $\eta \leq \min(n, m)$ , then

$$\frac{dZ}{dt} + \eta Z \leq \frac{K(r + \eta)^2}{4r},$$

where the quantity on the right hand side represent the maximum value attained by the parabola

$$\eta F + rF \left(1 - \frac{F}{K}\right).$$

Finally, the right-hand side of the above expression is thus bounded by a suitable constant  $\ell > 0$ , so that

$$\frac{dZ}{dt} + \eta Z \leq \ell. \quad (3)$$

The theory of differential inequalities [1] then ensured that

$$0 < Z(L(t), S(t), F(t)) < \frac{\ell}{\eta}(1 - e^{-\eta t}) + (L(0), S(0), F(0))e^{-\eta t}. \quad (4)$$

Thus as  $t \rightarrow \infty$ , we have  $0 < Z < \frac{\ell}{\eta}$ . Since the total population is bounded, also all the individual populations, solutions of (1), are bounded in  $\mathbf{R}_{0,+}^3$ .

### 3 Long term behavior of the model

In order to understand how the trajectories of the system behave in the long run, we study the equilibria of the dynamical system (1). There are only three equilibrium points. One is trivially given by the origin  $E_0(0, 0, 0)$ . Secondly we find the axial equilibrium point  $E_1(0, 0, K)$  in which only fish survive. The third one is the interior coexistence equilibrium point  $E^*$  with population levels given by

$$L^* = \frac{rm\Delta}{K}, \quad S^* = \frac{rg\Delta}{K}, \quad F^* = \frac{nm}{ag - bm}, \quad \Delta = \frac{(ag - bm)K - nm}{(ag - bm)^2} = \frac{K}{ag - bm} - \frac{nm}{(ag - bm)^2}. \quad (5)$$

The trivial and the axial equilibrium points always exist. The feasibility of the interior equilibrium point depends instead on the following conditions on the model parameters, namely

$$ag > bm + \frac{nm}{K}. \quad (6)$$

### 4 Stability analysis of the equilibria

The Jacobian matrix for the system (1) at an arbitrary point of the  $L, S, F$  phase space is given by

$$J \equiv \begin{pmatrix} -n - bF & aS & aS - bL \\ g & -m & 0 \\ bF & -aF & \frac{r(K-2F)}{K} + bL - aS \end{pmatrix}. \quad (7)$$

Our first result shows that for the system (1), total extinction is not possible. Indeed one of the eigenvalues of the Jacobian (7) at the origin is  $r > 0$  and so the origin is always an unstable equilibrium.

The eigenvalues of the Jacobian (7) at  $E_1(0, 0, K)$  are  $-r$ ,  $-(n + bK)$ ,  $-m$ . Hence, the axial equilibrium point  $E_1$  is always stable.

The Jacobian matrix (7) evaluated at  $E^*$  simplifies a little to give

$$J \equiv \begin{pmatrix} -n - bF^* & aS^* & aS^* - bL^* \\ g & -m & 0 \\ bF^* & -aF^* & \frac{-rF^*}{K} \end{pmatrix}. \quad (8)$$

From the latter, the characteristic equation of the matrix (8) can be obtained. It is given by the cubic

$$Q_1x^3 + Q_2x^2 + Q_3x + Q_4 = 0 \quad (9)$$

whose coefficients can be expressed in terms of the system parameters as follows

$$\begin{aligned} Q_1 &= K, & Q_2 &= nK + bF^*K + rF^* + mK, \\ Q_3 &= nrF^* - gaF^*K + bF^{*2}r + bF^*mK + mrF^* - bF^*KaS^* + b^2F^*KL^* + nmK, \\ Q_4 &= -gaF^{*2}r + ga^2F^*S^*K - gaF^*bL^*K + nmrF^* - bF^*KaS^*m + b^2F^*KL^*m + bF^{*2}mr. \end{aligned}$$

Due to the complexity of the above expressions, it is not possible to derive explicit stability conditions that can be easily interpreted biologically. Thus we turn to numerical simulations for our further analysis.

## 5 A particular case of system (1)

However, before proceeding further with the analysis of the model (1), it is interesting to consider the role on the dynamics of the system of other possible food sources for the fishes.

To emphasize the effect of alternate food sources, we assume here that the latter are not available for the growth of the fishes. In other words in absence of these other feeding possibilities, the net growth rate for fishes is negative, i.e if  $r = (f - e)$ , where  $e$  and  $f$  represent respectively the natural death and birth rates of the fishes, then  $f = 0$  making  $r = -e < 0$ . In such circumstance, we disregard also any intraspecific competition for the alternate food source. To make the corresponding term vanish, we could take the carrying capacity of the logistic growth to infinity, i.e  $K \rightarrow \infty$ .

With these assumptions, the model (1) reduces to

$$\begin{aligned}\frac{dL}{dt} &= -nL + aSF - bLF \\ \frac{dS}{dt} &= gL - mS \\ \frac{dF}{dt} &= -eF + bLF - aSF.\end{aligned}\tag{10}$$

The system (10) has only one feasible equilibrium point, i.e, the origin  $H_0(0,0,0)$ , which always exists.

The eigenvalues of the Jacobian of the system (10) at the origin are by  $-n, -m, -e$ . Thus the system (10) is always locally asymptotically stable around the origin. Since this is the only equilibrium of the system and the boundedness result of Section 2.1 continues to hold also in this case, the origin is also globally asymptotically stable. For this, observe that proceeding as in Section 2.1 we find  $\dot{Z} = (g - n)L - mS - eF \leq 0$ , so that as  $t \rightarrow \infty$ , we find  $Z \rightarrow 0^+$ , implying that each population in the system vanishes. In this case, we can thus consider  $Z$  as a Lyapunov function.

Thus we observe that in the absence of the alternative food source for the fishes, total extinction of the system is guaranteed. This however will never happen if the alternative food sources are available for the fishes.

## 6 Numerical experiments

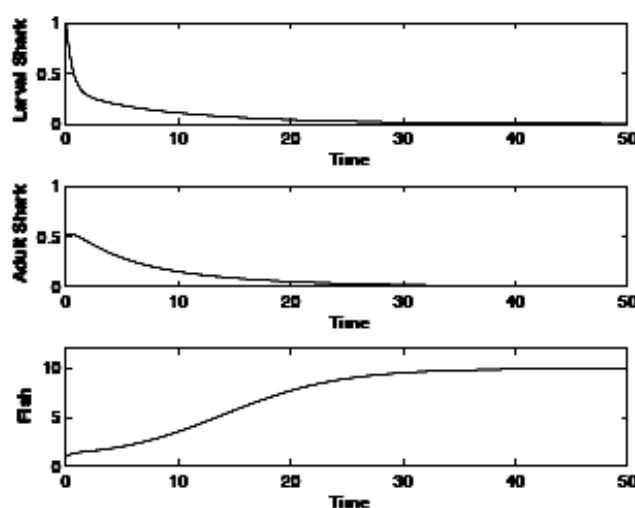
We performed extensive numerical simulation to substantiate our analytical results. To illustrate the results, we begin with the set of hypothetical parameter values given in Table 1. For these values,  $\Delta < 0$  follows, so that no interior equilibrium point exists. Thus the only stable point is  $E_1 = (0, 0, 10)$ , see Figure 1, which is also globally asymptotically stable.

Table 1: A set of hypothetical parameter values.

Parameters	values
$n$	0.7
$b$	0.8
$g$	0.25
$m$	0.3
$a$	0.7
$r$	0.2
$K$	10

We then vary the various parameters to study their effect on the dynamics of the system (1). Increasing the value of  $g$  to 0.38,  $\Delta$  becomes positive. In such case the interior equilibrium exists, namely the point  $E^* = (0.44, 0.56, 8.08)$  around which the system is locally stable, see Figure 2. A further increase in  $g$  to 0.45 increases the value of  $\Delta$ . This leads all populations to coexist through periodic oscillations, see Figure 3. Decreasing the value of  $m$  to 0.2,  $\Delta$  becomes positive. The interior equilibrium point becomes  $E^* = (0.1, 0.22, 9.3)$ . It is found to be locally asymptotically stable, see Figure 4. Further decrease of  $m$  to the value 0.1 increases the value of  $\Delta$ . The system now shows sustained limit cycles, see Figure 5. The same phenomenon is observed for changes in some other parameters too, which we do not report here.

Summarizing our findings, we find that for the parameter values making  $\Delta < 0$  the system is stable around the only equilibrium  $E_1$ , while no coexistence is possible, since  $E^*$  is not feasible. Instead, for small positive value of  $\Delta$ , the system is stable around the interior equilibrium point  $E^*$ . Finally for large positive values of  $\Delta$ , the populations of the system all coexist via periodic oscillations.


Figure 1: The system tends to the stable equilibrium steady state  $E_1$ .

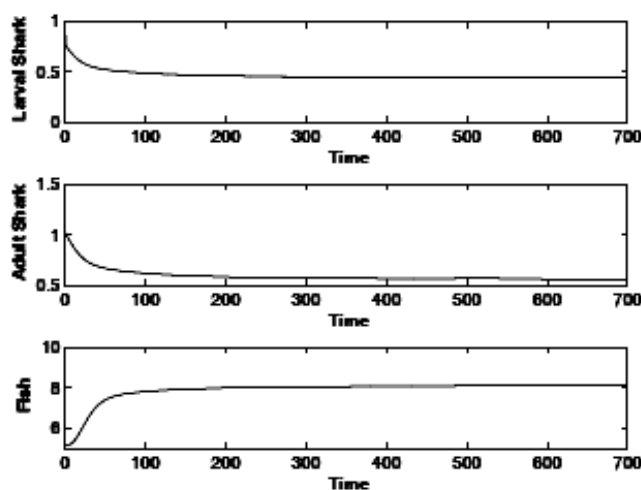


Figure 2: The system tends to the stable equilibrium coexistence steady state  $E^*$  for  $g = 0.38$

## 7 Interpretation

Our analytical and numerical results show that the stability of the system depends upon the parameter combinations leading to the value of  $\Delta$ . We found that the equilibrium with no sharks, either adult or at the larval stage, namely  $E_1(0, 0, K)$ , always exists and it is locally asymptotically stable. If the value of  $\Delta$  is negative then this fish-only equilibrium  $E_1(0, 0, K)$  is the sole stable equilibrium point of the system (1). But, if the value of  $\Delta$  is positive, then the interior coexistence equilibrium point  $E^*$  also is feasible. In this case for small values of  $\Delta$  then the system (1) is locally asymptotically stable near  $E^*$ . If  $\Delta$  attains larger values, the model (1) shows periodic oscillations. Thus, to maintain the stability of the system around the coexistence steady state the parameters have to be controlled so as to obtain a small positive value for  $\Delta$ . In particular we remark here that among the parameters appearing in the definition of  $\Delta$ , we find the mortalities of both larvae and sharks, which therefore have a great impact on the system's behavior.

Another important result we can gather from our analytical study concerns the role of alternative food source for the fish population. We observe that in the absence of the alternative food source total extinction of the system (1) becomes possible. This however will never happen in presence of alternative food sources for the fish population. Thus alternative food sources play a vital role in ensuring coexistence of the aquatic ecosystem.

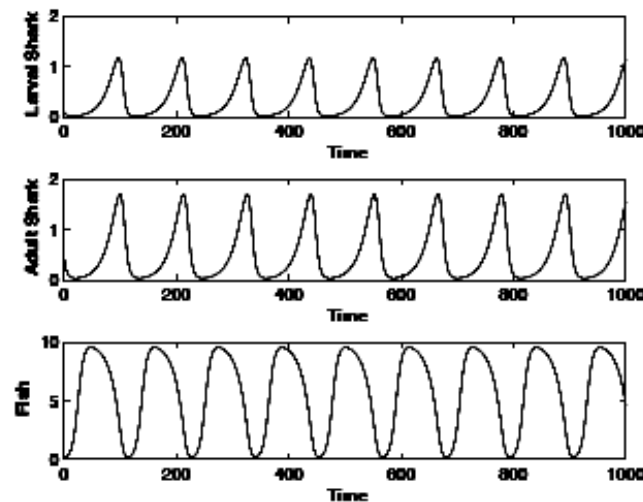


Figure 3: The system shows periodic oscillations for  $g = 0.45$ .

## 8 Some remarks on selective harvesting

We add here some considerations on the role of fisheries. Let  $H$  denote the harvesting function. Model (1) gets then modified as follows

$$\begin{aligned}\frac{dL}{dt} &= -nL + aSF - bLF \\ \frac{dS}{dt} &= gL - mS - H(S, F) \\ \frac{dF}{dt} &= rF \left(1 - \frac{F}{K}\right) + bLF - aSF - H(S, F).\end{aligned}\tag{11}$$

This situation corresponds to fishing made with boats and nets which capture reasonably sized fishes. Thus both  $S$  and  $F$  will be hunted. To model the fishery activities, let us make a very simple assumption, namely that  $H$  is a linear function of its arguments, so that in the second above equation  $H(S, F) = hS$  and in the third one,  $H(S, F) = hF$ . By suitably collecting terms then it is seen that the model (11) can be assimilated to the original model (1) in which the sharks' mortality  $m$  becomes larger, namely  $m + h$ , and the original net birth rate  $r$  decreases to the value  $r - h$ . Therefore, this type of fishing entails consequences that can be determined via the simulations of Section 6. Namely observe that with larger mortalities, the first term in the definition of  $\Delta$  increases, while the second one being a ratio of two quadratic terms does not change much. Hence  $\Delta$  increases and so do  $L^*$  and  $S^*$ . From (5) also  $F^*$  increases, since its numerator is quadratic. So apparently fishing helps the ecosystem to thrive. But as soon as  $m$  increases past the threshold,  $m^* \equiv ag\frac{1}{b}$ , the first term of  $\Delta$  becomes negative, and the coexistence equilibrium becomes infeasible, so that the system will suddenly go to extinction.

We consider now a second type of harvesting, namely the fishing of larvae. The model (1)



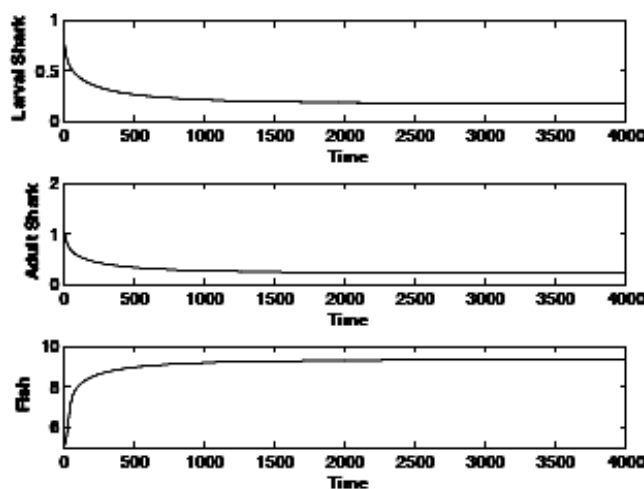


Figure 4: The system tends to the stable equilibrium coexistence steady state  $E^*$  for  $m = 0.2$

now becomes

$$\begin{aligned} \frac{dL}{dt} &= -nL + aSF - bLF - H(L) \\ \frac{dS}{dt} &= gL - mS \\ \frac{dF}{dt} &= rF \left(1 - \frac{F}{K}\right) + bLF - aSF. \end{aligned} \quad (12)$$

Assuming once again a linear harvesting function,  $H(L) = hL$ , fishing the larvae amounts to making their removal rate  $n$  larger,  $n + h$ . The long time effects of such situation can therefore once again deducted via the simulations of Section 6. In this case a larger  $n$  makes the second term of  $\Delta$  larger, so that this quantity becomes negative, past the threshold  $n^* \equiv K(ag - bm)\frac{1}{m}$ , making once again  $E^*$  infeasible and leading to the collapse of the ecosystem. Although in our formulation we spoke about sharks and their larvae, these results can be interpreted on fishing of the small fishes, which is nowadays regulated by law and limited to just a month in the spring, at least in the Mediterranean. The long term consequences of an indiscriminated fishing in such situations can be investigated via suitable simulations. The latter can therefore provide a scientific and rational support for the normative statements on this matter. However, to have a complete picture of the situation, also harvesting under other assumptions is needed, such as the more commonly used diminishing return functions, see chapters 1 and 2 of [4], and postponed to another investigation.

**Acknowledgement** Research supported by: MIUR Bando per borse a favore di giovani ricercatori indiani.

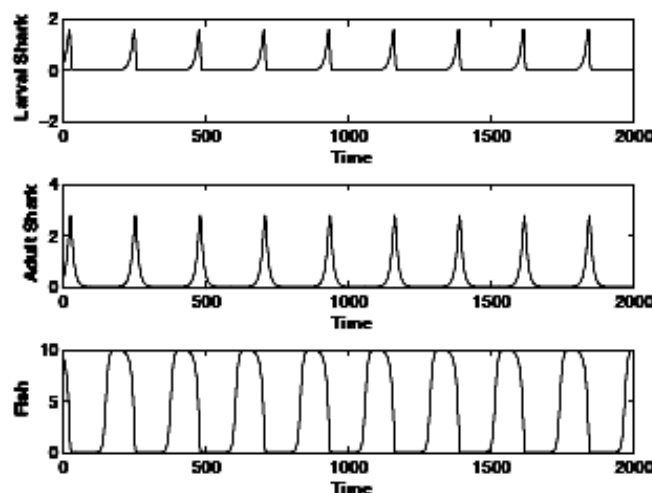


Figure 5: The system exhibits periodic oscillations for  $m = 0.1$ .

## References

- [1] BIRKHOFF, G., ROTA, G.C., Ordinary Differential Equations: Ginn Boston, 1982.
- [2] CARPENTER, S. R., J. F. KITCHELL, HODGSON, J. R., Cascading trophic interactions and lake productivity, *BioScience*, 35, 634–639, 1985.
- [3] CHASE, J. M., Abiotic controls of trophic cascades in a simple grassland food chain. *Oikos*, 77, 495–506, 1996.
- [4] CLARK, C. W., Mathematical bioeconomics: the optimal management of renewable resources, Wiley, New York, 1976.
- [5] HSU, S., B., HWANG, T., W., Global stability for a class of predator prey system, *Siam. J. Appl. Math.* 55, 763-783, 1995.
- [6] LEIBOLD, M. A., Resource edibility and the effects of predators and productivity on the outcome of trophic interactions. *American Naturalist*, 134, 922–949, 1989.
- [7] LEIBOLD, M. A., A graphical model of keystone predators in food webs: trophic regulation of abundance, incidence and diversity patterns in communities, *American Naturalist*, 147, 784–812, 1996.
- [8] LEIBOLD, M. A., J. M. CHASE, J. B. SHURIN, and A. DOWNING, Species turnover and the regulation of trophic structure, *Annual Review of Ecology and Systematics*, 28, 467–494, 1997.
- [9] OKSANEN, L., FRETWELL, S. D., ARRÜDA, J., NİEMELA, P., Exploitation ecosystems along gradients of primary productivity, *American Naturalist*, 118, 240–261, 1981.
- [10] RENSHAW, E., Modelling biological populations in space and time, Cambridge University Press, Cambridge, UK, 1991.
- [11] SAEZ, E., GONZALES-OLIVARES, E., Dynamics of a predator-prey model, *SIAM J. Appl. Math.*, 59, 1867-1878, 1999.

**Current address**

**Samrat Chatterjee, Dr.**

Dipartimento di Matematica, Università di Torino, via Carlo Alberto 10, 10123 Torino, Italy,  
tel. +39-011-670-3449 samrat\_ct@rediffmail.com

**Ezio Venturino, Professor**

Dipartimento di Matematica, Università di Torino, via Carlo Alberto 10, 10123 Torino, Italy,  
tel. +39-011-670-3449 ezio.venturino@unito.it



## A NEW ACTIVE-SET METHOD FOR LINEAR PROGRAMMING BASED ON TRANSFORMATION OF FEASIBLE DIRECTION ALGORITHM INTO UNCONSTRAINED MINIMIZATION PROBLEM

JURÍK Tomáš, (SK)

**Abstract.** A new active set algorithm for solving general linear programming problems is proposed. Its feasible search direction is transformed into the unconstrained minimization problem using the theory of duality. This transformation enable us to employ some simple and very efficient subgradient methods for nonlinear optimization. This transformation is employed only in the case that the projection of the objective value vector to some affine space is not a feasible direction. In this sense, the presented algorithm is an extension of a projective method [4]. The behavior and efficiency of the algorithm is demonstrated by the numerical experiments on the randomly generated problems.

**Key words and phrases.** linear programming, active set methods, duality, unconstrained optimization

*Mathematics Subject Classification.* Primary 90C05, 90C46; Secondary 90C59.

### 1 Introduction

The optimization is a natural process: it attracted both theoreticians and practitioners of all the world since the beginning. The linear programming (LP) is the easiest form among the huge class of the optimization problems. Many fresh ideas has been tried from its origination (see [1]). Besides the well known simplex and interior point methods, in the last decades many different sophisticated algorithms have subscribed to the wide spread of LP algorithms. The most popular nonstandard algorithms belong to the active-set methods: the aim is to find the constraints which are active at the optimal solution of the LP problem. Usually their iteration

points lies on the boundary of the feasible region and they find the feasible direction only using the active set of the present solution.

The two path algorithm [2] combines the solutions of the primal and the dual problem to speed up the convergence to the optimal solution. The authors showed that their algorithm is more than 100 times faster than the simplex method on the randomly generated problems. The pure active-set method Sagitta is based on the global viewpoint to the LP problem. This method is based on the observation that an intersection of the hyperplanes which contains the most-obtuse angles to the objective direction is an accurate candidate to the optimal solution. Obviously, this is true only in the special cases, but the method exploits the duality theory and Farkas lemma to gain the optimal solution quickly. Their successful computational results are remarked in [3].

The newest non simplex-type active set method has been proposed in [4]. Its idea is to find the feasible direction as the projection of the objective value vector into the space determined by the active constraints. In the case that this projection is not a feasible direction, the different feasible direction is obtained as a linear combination of an (unfeasible) projection vector and some vector which is perpendicular to the objective value vector. The encouraging results on randomly generated problems were also presented.

The algorithm proposed in this paper significantly improve the previous algorithm. Finding the feasible direction is transformed into an unconstrained nonlinear minimization problem that can be solved by the standard subgradient methods very efficiently.

## **2 The new method**

The main result is to present the deterministic algorithm for solving general linear programming (LP) problems. In the first subsection we propose some basic facts and a notation used, after that we describe the transformation into the unconstrained optimization problem. At the end of the main section we present the early numerical results we obtain on the randomly generated problems.

### **2.1 Preliminaries**

We consider the LP problem in the standard form (constraints are equalities and all variables are nonnegative)

$$\begin{aligned} c^T x &\rightarrow \max \\ Ax &= b \\ x &\geq 0 \end{aligned} \tag{1}$$

where  $0 \neq c \in \mathbb{R}^n$ ,  $b \in \mathbb{R}^m$  are the given vectors and  $A$  is the given  $m \times n$  matrix. In the rest of the paper we assume that the given problem (1) has an optimal solution and some feasible solution  $x^0$  is available. If this is not the case, we can solve some additional LP problems similarly to the two-phase simplex algorithm. (To detect the unboundedness of a feasible problem it is sufficient to show that the dual problem is infeasible.)

The lower index  $i$  of the vector  $v$  denote its  $i$ -th element, while the upper index usually denotes the order of iteration points. The proposed algorithm starts from a given feasible point  $x^0$  such that  $Ax^0 = b$  and  $x^0 \geq 0$  and it generates the sequence  $\{x^k\}_{k \in \mathbb{N}}$  of some feasible points which converges to the optimal solution  $x^*$  of the problem (1). The set of all indices

$$N^k = \{i \in \{1, 2, \dots, n\} : x_i^k = 0\}. \quad (2)$$

is called an active set  $N^k$  corresponding to the feasible solution  $x^k$  (We note that here we omit the equality constraints  $Ax = b$  of the problem (1)). The feasible direction  $v^k$  that force the current feasible solution  $x^k$  to the next iteration point  $x^{k+1}$  has to be the projection of the objective value vector  $c$  onto the set

$$Av = 0, \quad v_{N^k} \geq 0. \quad (3)$$

The easiest way how to calculate the vector  $v : \|v\| \leq 1$  that fulfills (3) and maximizes  $c^T v$  is to set  $v$  as a projection vector  $c$  to all active constraints, e.g.,

$$v^k = \left( I(n) - A^{kT} (A^k A^{kT})^{-1} A^k \right) c, \quad \text{where} \quad A^k = \begin{pmatrix} A \\ I_{N^k} \end{pmatrix}. \quad (4)$$

Here,  $I(n)$  denotes the identity matrix of order  $n$  and  $I_{N^k}$  denotes the rows of the matrix  $I(n)$  which indices lies in the set  $N^k$  in that order. The drawback of the formula (4) is that it is a correct definition if and only if the matrix  $A^k$  is regular. Otherwise we have to calculate the (force) vector  $v^k$  in a different manner.

## 2.2 Unconstrained optimization problem

Our aim in each iteration is to find a vector  $v$  such that  $\|v\| \leq 1$  (for simplicity the subscript  $k$  will be omitted) which fulfills (3) and  $c^T v > 0$ . Then we have  $c^T x^{k+1} > c^T x^k$ , where  $x^{k+1} = x^k + \lambda v$  for some  $\lambda > 0$ .

In this subsection we show how to reformulate this problem into some unconstrained optimization problem. Our reference book for LP theory is [5]. The trick is to relax the condition  $\|v\| \leq 1$  and replace it by  $\|v\|_1 \leq 1$ , e.g.,  $-1 \leq v \leq 1$ . Then the primal-dual pair of the modified problem for a vector  $v$  has a form ( $N = N^k$ ,  $B = \{1, 2, \dots, n\} \setminus N$ )

$$\begin{array}{ll} c^T v \rightarrow \max & x + y \rightarrow \min \\ Av = 0 & (A^T w - x + y)_N \geq c_N \\ v_N \geq 0 & (A^T w - x + y)_B = c_B \\ -1 \leq v \leq 1 & x, y \geq 0 \end{array} \quad (5)$$

where  $x, y \in \mathbb{R}^n$  are additional variables corresponding to the constraints  $-1 \leq v$  and  $v \leq 1$  respectively. One can easily see, that the dual problem (5) is equivalent to the unconstrained minimization problem

$$f_N(w) = \sum_{i \in N} \max\{0, (c - A^T w)_i\} + \sum_{i \in B} |(c - A^T w)_i| \rightarrow \min.$$

If we want to apply some standard optimization routines for solving an unconstrained optimization problem (i.e. subgradient methods) we have to smoothen the function  $f_N(w)$  (to make it differentiable). The simple perturbation yields to the function

$$f_{\mu,N}(w) = \frac{1}{2} \sum_{i \in N} \left( (c - A^T w)_i + \sqrt{\mu + (c - A^T w)_i^2} \right) + \sum_{i \in B} \sqrt{(\mu + (c - A^T w)_i^2)} \quad (6)$$

for sufficiently small positive parameter  $\mu$ . This problem can be efficiently solved by the subgradient methods (see [6]).

After we calculated the optimal solution  $w$  of the problem (6) which is the reasonably precise approximation of the dual problem (5) we employ the complementary slackness conditions (c.f. [5]) to find the corresponding primal optimal solution  $v$ . The remaining variables  $x, y$  of dual problem (5) can be calculated from the residual  $r = c - A^T w$  as

$$x_i = \begin{cases} 0, & \text{for } i \in N \text{ and } r_i \geq 0, \\ 0, & \text{for } i \in N \text{ and } r_i \leq 0, \\ 0, & \text{for } i \in B \text{ and } r_i \geq 0, \\ -r_i, & \text{for } i \in B \text{ and } r_i \leq 0, \end{cases} \quad y_i = \begin{cases} r_i, & \text{for } i \in N \text{ and } r_i \geq 0, \\ 0, & \text{for } i \in N \text{ and } r_i \leq 0, \\ r_i, & \text{for } i \in B \text{ and } r_i \geq 0, \\ 0, & \text{for } i \in B \text{ and } r_i \leq 0. \end{cases}$$

In our case, from the complementary slackness conditions we can easily derive that the optimal solution  $v$  is defined by the rule

$$v_i = \begin{cases} 1, & \text{if } r_i > 0, \\ -1, & \text{for } i \in B \text{ and } r_i < 0, \\ t_i, & \text{for } i \in N \text{ and } r_i \geq 0, \end{cases} \quad (7)$$

where  $t$  is a solution of a regular linear system such that  $Av = 0$ .

### 2.3 Algorithm

The aforementioned ideas leads to this new algorithm for solving nondegenerate linear programming problems. Given a feasible solution  $x^0$  of the problem  $\max\{c^T x : Ax = b, x \geq 0\}$  we initialize our algorithm:



```

while (true)
    construct the sets  $N = N^k$  and  $B$  defined by (2)
    calculate projection  $v^k$  defined by (4)
    if ( $v^k = 0$ )
        solve unconstrained optimization problem (6)
        calculate corresponding primal feasible direction  $v^k$  due to (7)
    endif
    if ( $v^k = 0$ )
        break (optimality reached in  $x^k$ )
    endif
     $x^{k+1} := x^x + \lambda v^k, \lambda = \min \left\{ \frac{x_i^k}{-v_i^k} : i \in \{1, \dots, n\}, v_i^k < 0 \right\}$ 
     $k = k + 1$ 
endwhile

```

## 2.4 Numerical results on randomly generated problems

The proposed algorithm has been experimentally tested to claim its efficiency. We describe our numerical experiments and present computational results which demonstrate the efficiency of the new algorithm on randomly generated linear programs. We confine our experiments only to the dense random problems that were small and medium in size. The calculations have been made in the environment Matlab 6 performed on a PC with 2 GHz Intel Core<sup>®</sup>2Duo processor, 3 GB of RAM and the Windows XP operating system. The generated problems were in the form

$$\max \{c^T x : Ax = b, x \geq 0\},$$

where  $A \in \mathbb{R}^{m \times n}$ ,  $c, x \in \mathbb{R}^n$  and  $b \in \mathbb{R}^m$ . The matrix  $A = (a_{ij})_{m \times n}$  and the vector  $c \in \mathbb{R}^n$  were dense with  $a_{ij} \in [-0.3, 0.7]$ ,  $c \in [-0.3, 0.7]$ . In order to facilitate the computations, a vector  $b$  has been set to  $Ae$ , where  $e$  is all-one vector of size  $n$ . Ten problems were generated for each problem size and the average values are presented. The feasibility and precision tolerance were set to  $10^{-10}$ . The new algorithm was compared to the implemented procedure *linprog* in the environment Matlab. Since the presented algorithm belongs to the active set method, we used the procedure *linprog* with 'LargeScale' option turned off (simplex method).

The first two columns represent the size of the problems - number of constraints and number of variables respectively, the next columns represent the minimum, maximum and the average values for the CPU time (in seconds) and the number of iterations for both compared methods. The last column indicates the average number of solving the unconstrained minimization problem. From these results can be seen that the calculation time is comparable to the implemented Matlab *linprog* function. With the higher dimension, this ratio is going better for the new algorithm and this encouraging empirical observation that can stimulate the further improvements of the described algorithm.

		linprog						new algorithm						
size		time			iter			time			iter			unc.
$m$	$n$	min	max	avg	min	max	avg	min	max	avg	min	max	avg	avg
10	100	0.18	0.28	0.20	100	90	97	0.12	0.79	0.29	91	102	93	2.1
20	100	0.14	0.22	0.16	80	88	83	0.20	0.81	0.49	81	90	84	2.7
30	100	0.13	0.25	0.14	70	82	74	0.25	1.00	0.63	71	77	74	2.8
20	200	0.94	1.02	0.96	183	189	184	0.94	1.67	1.16	184	190	186	2.9
30	200	0.81	0.98	0.85	171	182	176	1.12	2.40	1.56	172	183	177	3.3
40	200	0.75	0.84	0.78	164	172	167	1.64	3.00	2.34	167	177	171	3.7
30	300	2.69	3.28	2.83	271	288	277	2.50	3.54	3.08	273	284	278	3.4
40	300	2.56	2.78	2.66	262	275	268	3.23	6.12	4.28	263	283	270	4.1
30	400	6.23	7.07	6.67	372	386	377	4.48	7.07	5.34	375	391	380	3.8
40	400	5.97	6.28	6.17	362	374	368	5.84	7.89	6.97	363	376	369	3.9

Table 1: Computational results for solving randomly generated LP problems

### 3 Conclusions

This paper fully describes the new idea of transformation of the feasible search direction, which is the most difficult (and therefore the most interesting) part of all active set methods. Our transformations is just slightly relaxation of the projection of the objective value function. The duality theory and the complementary slackness conditions help us to create a brand new algorithm for solving nondegenerated linear programming problems. A complete theory for the new algorithm has not been presented here, its convergence is still an open question. However, the early calculations provide a sufficient motivation for further research.

### Acknowledgment

The work of the author was supported by the grant No. UK/410/2008.

### References

- [1] SLOAN, S., W.: *A steepest edge active set algorithm for solving sparse linear programming problems*. International Journal for Numerical Methods in Engineering, Vol. 26, pp. 2671–2685, 1988.
- [2] PAPPARIZOS, K., SAMARAS, N., STEPHANIDES, G.: *A new efficient primal dual simplex algorithm*. Computers & Operations Research, Vol. 30, pp. 1383-1399, 2003.
- [3] PALOMO, A., S.: *The sagitta method for solving linear programs*. European Journal of Operational Research, Vol. 157, pp. 527–539, 2004.
- [4] JURÍK, T.: *Computational experience with a new active-set method for linear programming*. Journal of Electrical Engineering, Vol. 58, No. 7/s, pp. 24–27, 2007.
- [5] VANDERBEI, R.: *Linear Programming: Foundations and Extensions*. Second Edition, Springer, New York, 2001.

- [6] BERTSEKAS, D. P.: *Nonlinear Programming*. Second Edition, Athena Scientific, 1999.

**Current address**

**Tomáš Jurík, RNDr.**

Department of Mathematical Analysis and Numerical Mathematics,

Comenius University,

Mlynská dolina, 842 48 Bratislava, Slovak Republic

phone: +421 2 602 95 403

e-mail: jurik@fmph.uniba.sk



## OBJECT - ORIENTED PROGRAMMING LANGUAGES AS TOOLS FOR FORMULATIONS OF SYSTEM ABSTRACTION

KINDLER Eugene, (CZ), KŘIVÝ Ivan, (CZ)

**Abstract.** The classes that the object-oriented programming languages allow to define correspond to exact abstract notions. The instances of classes correspond to exactly behaving entities. The so called life rules that some of those languages allow including as integral components of classes correspond to algorithms that run in parallel. The life rules can be related to a certain abstraction of Newtonian time, which allows formulating exact description of dynamic systems, i.e. models of entities that are studied by natural and/or technical sciences and even by social sciences and humanities. The sets of classes can represent exact theories. In case such a language is also block-oriented, the exact theories can be formulated and handled as concepts, too. When the applied language admits the life rules, a certain dynamic development of such a theory can be exactly formulated. Moreover, the synthesis of the block orientation with the object orientation allows describing dynamic systems that handle exact theories. If the applied object-oriented language is consistent and independent of computer essence, it represents a true mathematical language, able to be used for describing very complex systems (including intelligent ones).

**Key words.** Object-oriented programming, exact theories, nesting theories, complex systems

*Mathematics Subject Classification:* Primary 68N15, 68N19, 68N30, 68U01, 68U20; Secondary 68T35, 93B07.

### 1 Introduction – Development of the Languages of Constructive Mathematics

Beside the languages of formulas broadly applied by the mathematicians, the constructions were in focus of mathematics and their authors tried to design an formal language that could be used for exact and simply decipherable describing of such constructions. Already the pupils of elementary schools learned a language for declaring Euclidean constructions. Naturally, that language was very limited. Nevertheless, it showed a certain way to the further development. Note that when phases of the described construction had to be repeated, simple phrases of a natural language were applied. Although the natural languages are evidently non-exact tools, the use of the

phrases in the description of the constructions was rather limited, when compared with the use of natural language in the description of definitions, theorems and their proofs occurring in conventional mathematics. The other formal languages used for describing constructions were limited, too (e.g. those for solving equations or for analytic computing definite integrals – the forms oriented for computation at electromechanical desk top machines can be included). Efforts to formulate a language that would be much more general (even universal), i.e. that would be able to be used for the description of an essentially larger set of constructions led to inventing abstract tools, the commonness of which was paid by illegibility in case of practical applications (sometimes rather popular in Czechoslovakia Markovian algorithms [1], lambda calculus, Turing machines, general recursive functions etc.).

Since the end of the World War II, the efforts started from another source, namely from programming digital computers. Their programs were representations of true constructions but at the first phase, when programming was manually performed in machine code or in similar “language” like that of symbolic addresses, the applied tools were as illegible as those mentioned above. Nevertheless, the debugging or programs profiting of the physical reality of computers, stimulated the further development and unlatched narrow horizons that limited the thinking of those using the general tools for construction describing, which had arisen independently of computers.

Automatic programming soon discovered languages that seemed to be rather universal and more legible, like Flowmatic or Mathmatic, introduced to old Univac computers in the mid of the sixties of the last century. The main step in the development consisted in a synthesis of arithmetic formulas with simple words and phrases of the English natural language. That development continued further (let us remind Fortran II. and IV. and Algol 58) and nowadays one can say that it was crowned by Algol 60 [2].

## 2 From Algorithms to Processes and Concepts

The described development led to the possibility to formulate the definition of algorithms working with predefined “standard” entities, like numerical, text or Boolean ones. However, the old language used for Euclidean constructions concerned geometrical entities, i.e. something that was rather different from standard ones. Some authors tried to overcome that unpleasant situation in a certain general sense, namely by introducing data structures. In general, data structure is composed of standard entities and possible pointers to other data structures. So languages as COBOL, PL/1, ALGOL W, ALGOL 68 or early versions of PASCAL were designed and implemented. The general idea rooted in the illusion that an exact abstraction from any concept should be a set of standard entities and a priori formulated possible relations.

But almost contemporarily to the development the just mentioned idea and illusion, another development existed, almost neglected by the scholars oriented to programming. This development took into account that the entities often behave dynamically in a certain autonomous manner (in a metaphor: they live) – such entities are often understood as instances of a concept for which the autonomous behavior (“life”) is a proper component of its semantic contents. Most probably, the mentioned neglecting rooted in the fact that such entities use to be rather distant from mathematical ones, having been studied by sciences other than mathematics and physic. The essential stimulus for the development was discrete event simulation. In simulated systems parallel processes exist, which have to be mapped as algorithms in the corresponding simulation models. But a projecting of such parallel processes to a conventional (monoprocessor) computer leads to a very sophisticated algorithm and, therefore, the simulationists invented so called *process oriented discrete event*

*simulation languages* in order to facilitate the fabrication of such algorithms. The user of such a language does not need describe what the computer should perform, but he describes the simulated system as a dynamic structure of elements that “live” and the description is automatically converted into a program that can run at the concerned computer. The first language of that sort was GPSS [3]. It offered describing systems as composed of entities that could enter and leave the system during its existence and that are structures of data (so called *attributes*) and of “*life rules*” that had form of algorithms interleaved with so called *scheduling statements*, expressing duration of particular phases of the “life” and – indirectly – switching among algorithms performed by different entities. Rather poor abilities of GPSS provoked further development of the process oriented discrete event simulation languages, passing over SOL [4], SLANG [5] etc. to language called initially SIMULA and later SIMULA I [6, 7] in order to be distinguished from a much different and later developed object-oriented language called SIMULA 67. It is to observe that the title of [6] characterizes SIMULA as a language for the description of systems and not for programming.

Note that not every discrete event simulation languages were process oriented. Among popular examples, both main versions of SIMSCRIPT and all versions of GASP can be presented.

The result of this development consisted in formal languages that offered exact description of clusters of concepts concerning in a common target – reflecting a structure and dynamics of a certain set of discrete dynamic systems.

### 3 SIMULA 67

In 1967 the first official presentation of a new language SIMULA 67 was presented [8]. This language has several essential properties, discussed in the following text.

(1) It introduced the relation class-subclass, called often *subclassing* or *specialization*, which reflected the relation of enriching the content of a concept or – inversely – the reduction of the set of concept instances. Such a relation is widely used, when a new concept is introduced (“concept *A* is defined as concept *B* for that so and so new properties can be observed”). It is possible to state that without that relation the human society (including science, communication and control) could not exist.

(2) It offered a possibility to declare algorithms as components of contents of the classes, to give them names and to call them to work for any instance of the concerned class. Later, such algorithms were called *methods*.

(3) It introduced *dot notation* (called also remote identifying) for expressing attributes and methods: *A.s* represents either attribute *s* of instance *A* (in case *s* is an attribute) or “let *A* perform *s*” (in case *s* is a method). In the first case, the dot notation serves for expressing what the logicians call determination (question: “What *s*?”, answer “*s* of *A*”, i.e. a phrase that is often expressed “*A s*” in English and in some other languages). In the second case, the dot notation serves for expressing simple phrase (*A* be its subject and *s* its predicate). Moreover, the methods can have parameters (like procedures) and thus *A.s(p)* can serve for expressing phrase where *A* is subject, *s* predicate and *p* object or complement. That does not exclude other interpretations, where *s* can be e.g. a copula or a preposition.

(4) Methods can be introduced as *virtual*, i.e. so that they are considered as meaningful, supposing their proper content will be declared in subclasses. It can serve to introduce certain sort of adaptation to the context, because such a method can adapt its work according to the “sort” of the object standing in front of the dot (“sort of *A*” means the class, of which *A* is an instance).

(5) Life rules can be added to the declaration of any class and may be enriched in declaration of any of its subclasses. Calling methods and forming new class instances can occur in the life rules.

(6) Beside the virtual methods, also virtual tasks of the transition (called go to statements) among the components of the life rules are allowed.

(7) Declarations of classes can be anywhere when declarations of attributes and methods can occur. A class that contains such class declarations is called *main class* and can serve for introducing formal theories, offered for manipulation with more than one class and its instances.

(8) Full block oriented facilities introduced in ALGOL 60 exist in SIMULA 67. For a block, not any variable and subprogram can be introduced but also classes. It allows handling with local classes, or – viewed in another way – with concepts that are related to a certain context (block) while they can differ from concepts that have the same names.

The introduced properties were included also into a refined definition of SIMULA 67 [9]. This language was later presented as international standard [10] and its name was simplified to mere SIMULA (when SIMULA I was generally forgotten because its users decided to use SIMULA 67).

## 4 Object-Oriented Programming

In the 80-ies of the 20th century, the properties (1) – (4) were taken as characteristics of the *object-oriented programming*. But already in 1968 one of the authors of SIMULA 67 called the mathematicians' attention to the large horizons that the programming languages disposed by the mentioned properties for exact representation of concepts [11].

He demonstrated the ideas at SIMULA 67, as in 1968 no other language with properties (1) – (4) existed. Nowadays, there are many object-oriented languages and among them popular C++, JAVA and newer versions of PASCAL, but it is difficult to consider them as suitable tools for the formulation of concepts. The reason is that while SIMULA 67 is completely independent of the computer at that it is used, the other languages are bound with it. So they are inconsistent and rather distant from being tools for the representation of exact concepts (distant from mathematics), on the other hand being nearer to computer science. Let us illustrate the “non-mathematic” properties of the most popular object-oriented languages.

PASCAL, even in its object-oriented versions, has no garbage collector, and so one can let liquidate an object to that still some pointers exist. After the liquidation of the object, the memory, if required, it needed is given to disposal for possible new objects. Suppose an object  $A$  have been liquidated but a pointer  $P$  to it remained; the segment  $M$  of the computer memory, at which  $A$  was represented, can be used for other purposes. Suppose another object  $B$  is then generated, belonging to a class different from that of  $A$ . And suppose (a part of)  $M$  is used for storing  $B$ . Then  $M$  is completely re-structured, the borders between items are in general displaced and when  $P$  demands some item of  $A$ , it obtains a nonsense.

Similarly as PASCAL, C++ has no garbage collector; thus for C++ the same illustration as for PASCAL holds. But moreover, C++ is intended as a language suitable for the design and implementation of operating systems and the consequence is that C++ has to be sometimes for disposal as an autocode. Truly, in case  $B$  is an object introduced in a software formulated in C++, one can work with the internal (computer) address  $\alpha$  used e.g. for the first item of the representation of  $B$ . In general, when  $B$  is an abstraction of something that exists (or may exist) independently of the computer at that  $B$  is hand,  $\alpha$  has no relation to  $B$ , moreover,  $\alpha$  can have different values for the same model running at different computers or even for the same model running at the same



computer at different jobs or tasks. Another example can give an opposite step – e.g. determining the object stored at a part of computer storage using address  $2\alpha$ .

JAVA is in a certain sense an enlargement of C++ and so the texts written in it are open to suffer similarly as that written in C++. Besides, JAVA offers other obstacles related to its bindings with the carrying computer hardware, namely concerning the parallelism and so called passports of the objects, i.e. with data structures that have only internal, computer-based interpretation related to an object (in other words: a passport of an object  $A$  tells nothing on the state and context of  $A$  relating to the world where  $A$  is supposed, but only something about the internal state of the computer model of  $A$ ).

The use of the mentioned three popular programming languages as tools for the representation of concepts is further disadvantaged by a number of dialects that are broadly applied. The same problem exists also for another popular object-oriented language SmallTalk that exists also in different variants. Moreover, that language is “self-defined” (defined by means of the concepts proper to its own semantics and according to that implemented), which enables its users that by making a programming error they damage the whole semantics (and therefore the applicability – both as a programming tool and as a tool for concept representation) of the language. At the present days, beside SIMULA 67 only BETA [12] seems to be a logically consistent tool for the representation of concepts, but its heavy drawback is a great distance of its syntax from both current communication conventions in mathematics and natural English.

The object-orientation of C++ and PASCAL is isolated and not synthesized with the process orientation. That causes further obstacles, as the “lives” that are to be modeled (often simulated) at a monoprocessor computing system and that are viewed as developing contemporaneously in the modeled universe, have to be broken into “events” and considered as carried by hypothetic objects that have no pattern in the modeled universe and that are considered as “living” a mere instance of zero duration”. Evidently, such languages are not suitable to be used for the concept formulation. Note that all simulation tools, the implementation of which is purely based at C++ or Pascal, are of the just mentioned sort. Some experiments trying to enlarge those languages to be processed oriented must have use of procedures programmed in machine code; beside other, such a manner evidently binds the mentioned languages with computer hardware, retreats them from their abstract function and draws them near to computer science as mere programming tools.

It is interesting that essential tools concerning the process orientation of JAVA are made in a similar way (by means of machine code or tools that are near to machine code). The basis, i.e. the system of threads, is essentially hardware oriented aspect and in case one would like to adapt it for representing and/or modeling parallel processes, he has to express a certain way just from hardware to this universe [13].

The consistent form of the object-oriented languages independent of computer hardware and the elegant manner of accumulation of large amount of exact information enabled by using such languages stimulated an idea to propose similar tools for “specification languages”, i. e. languages that had to serve only for abstract exact description of systems but on which one never hoped to be implemented. SIMULA directly served for defining language DELTA [14] (even one of the authors of SIMULA was among the authors of DELTA). Nevertheless, neither DELTA nor other proposed languages were successful. It is naturally better to describe a system using a language for that an implementation (a compiler) exists. Curious information was presented by another author of SIMULA in [15]. He described a system for postal package sorting with faults (and therefore many cycles) in a specification language and in SIMULA, then he applied the formal definition of the specification language semantics in order to get another way to the description in SIMULA. After analyzing it, he showed that the normal analysis of a system, which makes a human during

preparing the corresponding classes and applying them for writing the simulation model, is essentially different from the process of formal “translation”.

## 5 Logic Programming

Logic programming is not so broadly spread as object-oriented one but many of its aspects are close to formal logic. Therefore, relating to formulation of exact concepts, logic programming demands to be discussed beside the object-oriented one. From that point of view, languages facilitating logic programming (PROLOG) seem to serve as a base for formulation of exact concepts, but judged from many other points of view, logic programming is not suitable. Let that be explained in certain details:

(a) Although the language of formal logic (of first order predicate calculus) serves as a good exact base for introducing and growing mathematical theories axiomatized in conventional way (i.e. almost all mathematical theories that arose independently of computing technique), it appears heavy-handed and cumbersome for defining processes running over time.

(b) Exerting a great attention is necessary so that one should not forget some evident aspects of real processes based on abstractions of those existing in the real world. As an example, we can present the axiom stating that when an object enters a certain state it has to remain in it during some positive time interval. If that is not formulated then concepts like the following one is admissible: a human who is a patient in a given hospital during time interval  $(t_1, t_2)$ , where  $t_1 < t_2$ , may occupy one bed at time instants the value of which is a rational number, and another bed at time instants the value of which is an irrational number.

(c) The form of such a language allows producing incomplete definitions in the sense that a phase of a process (of life rules related to an object or to a class of objects) is formulated but description of its continuation is forgotten.

(d) The computer models based on logic programming run slowly because of a so called unification (ordinary connecting bound variables occurring at quantifiers belonging to different formulas).

(e) The present intellectual state of normally educated and thinking persons prefer algorithmization to axiomatization.

## 6 Handling of World Viewings

The object-oriented languages that are also process-oriented can be taken as much more suitable agent-oriented languages; their great suitability consists in their universality which – among other – allows defining other agent-oriented programming tools that are not so general and that are often tailored to some limited conception of agents. Naturally, when – like SIMULA – such a language is strictly separated from the expressing tools concerning the hardware at that the models written in the language can operate, it appears a suitable agent-oriented tool for representation of concepts, which is prepared to be used for representing concepts of entities behaving in time and – moreover – for representing dynamic objects whose “life” exists contemporaneously to other similar objects but cannot be measured as existing in (Newtonian) time (as the quasi-parallel sequencing arisen by the generalization of scheduling in time, being implemented in SIMULA and BETA).

Nevertheless, the fresh demands of the exact branches of science, technology and humanities, stimulated by vehement development of the computing technique, ask more. Nowadays, two sources of those demands seem to be identifiable:

(i) There are often different opinions on the same subject, where the expressed concepts are called in the same manner but the contents of the concepts can differ in different opinions. Such phenomenon can be observed in a large spectrum of the world viewings, beginning from simple views to the space (one agent uses e.g. Cartesian geometry while another one applied spherical one), going over technical statements like fuzziness/equivocality of some future events (an example: one supposes it can be expressed in Gaussian terms, another one prefers using of old histograms) and ending in complex hypotheses on financial crises, on global warming and energy sources. When a process of communication among entities with different opinions on the same phenomenon or project should be exactly described (or even modeled), then it is to respect that each of such entities has its private fund of concepts, while for describing the communication as a whole (or – possibly – for describing a system in that such a communication forms a component) another fund (called global) of concepts should be at disposal. Thus the “private” concept funds have to be seen as nested in the global fund.

(ii) The objects that come under one’s exact studies are more and more complex and tend to be equipped by information processing facilities – either by computing elements (computers) or by human thinking, having use of experience, logical derivation and imagining. The human makes those activities with support of more or less general notions and the merit of the object-oriented programming paradigm stimulated the persons who equip the computing elements of the systems to have use of this paradigm. So in both the cases, the exact representation of concepts “private” for the thinking humans and/or for the computing physical elements comes out as a suitable and efficient technique of the professional manner of studying, modeling and communication. Such thinking and/or computing element has a fund of its “private” concepts. The funds should be nested inside another fund of concepts, serving for formulation of the whole system in that the computing and/or thinking elements exist as its components. That last fund corresponds to the global fund introduced in (i) and so it will be called in the further consideration.

The phenomenon of nesting of private funds in the global ones can be a bit refined and concretized.

Firstly, any private fund  $P$  is owned by a certain element  $E$  and thus it is suitable (and near to the reality from that the formulation is abstracted) to consider  $P$  as an attribute of  $E$ . In general, more elements like  $E$  can be in the described system  $S$  and so it is reasonable supposing  $E$  an instance of a general concept  $C$  of “objects similar to  $E$ ” (such objects can differ by some attributes, abilities or even – as it will be shown later – by some fine details of the private concepts used by them).

Secondly, both the global concepts and the private ones can be formulated as classes respecting the paradigm of the object-oriented programming and possibly facilitated by life rules. The funds of concepts can be mapped as sets of classes, but some languages (like SIMULA and BETA) offer so called *main classes*, i.e. classes, for which not only attributes, methods and life rules can be declared as their contents, but other classes, too (called nested ones). The main classes make possible expressing relations among the nested classes they contain and among them and their particular instances.

A simple example can be presented by main class *geometry*, meant as Euclidian plain geometry, which contains classes like that of points, that of lines, that of circles etc. Introducing class of circles, one can designate it as defined by its *radius* and *center*, where *center* is an instance of the class of points. Among the methods executable by the points, one could introduce the tests

like “a point is inside a circle”, “a point is outside a circle”, “a point is at the circumference of a circle etc.). Beside the classes nested in *geometry*, one can introduce attributes pointing to the origin of coordinates and to the coordinate axes and – in the life rules – generating of an instance of the class of points and two instances of the class of lines and assign them to the just mentioned attributes pointing to the coordinate system components. As the computation of real values is limited by the finite length of the computer word, it is difficult exactly to test e.g. whether two points are in geometrical sense equal, and so class geometry should be defined as having a real attribute, e.g. *epsilon*, serving for tests like a point is equal to another point means that the distance between both the points is less than *epsilon*. The presented example can be simply changed to that concerning a production system or a transport one in place of the Euclidean plane.

Thirdly, class *geometry* is an efficient stimulus for considering a main class as a formal theory. In case a main class  $G$  is carried as an attribute by an object it represents a situation that the object is a carrier of a certain sort of thinking that conforms to theory  $G$ . Let us consider a system  $S$  containing elements that carry theories (e.g. persons or machines that plan some moves and thus use class *geometry*). Then  $S$  can be described using a main class that may be interpreted as another theory  $T$ . The theories like  $G$  are then nested into  $T$  in the sense that  $T$  concerns (beside others) entities carrying by their own theories like  $G$ . If one quests for other scientific tools able to study theories of entities that carry other theories, only theoretical arithmetic of the natural numbers offers that tool with its technique of gödelisation. Evidently, such a way is so elementary that every idea on its practical application is an illusion. Some applications, namely for the simulation of anticipatory systems that use their own simulation models to anticipate possible future consequences of their decisions, were indicated in [16].

## 7 Blocks and Their Impact

Already ALGOL 60 was equipped by the whole apparatus of blocks [2]. The block is a part of algorithm that is able to manipulate with an entity that it inaccessible outside that part. In ALGOL 60, such an entity could be a procedure or a variable. Such entities were called entities *local* (more precisely: global in the block), while the other entities (accessible also outside the block) were called *global* (more precisely: global for the block). The first reason to introduce local entities was storage economy (the local entities could be removed from the operation memory when the algorithm operated outside the block), but soon other facilities of them were discovered:

( $\alpha$ ) A description of an algorithm can contain more than one block so that its run can be present at least in one of them. When the run is so in a block it is a certain (very poor) image of a “special phase of the algorithm’s real existence, during which it can use as special thinking tools the global entities”. The same name can be used for identifying a local entity in a block and for identifying a local entity in another block. Both the denoted entities are semantically different and that may be an image of a situation that the algorithm “respects different abstractions of a thing denoted by the name during each of the phases of thinking”.

( $\beta$ ) The operation component of a block is composed of statements; a block is viewed as a statement, too, and therefore a block  $A$  can contain another block  $B$  among its components. Then  $B$  is called a *subblock* of  $A$  and the entities local in  $A$  turn global for  $B$ . In case an entity local in  $A$  has a name used for an entity local in  $B$ , a so called *name conflict* takes place. For such a situation, ALGOL 60 introduced a rule, according which the name denotes the entity local in  $B$ ; the consequence followed, that the entity global for  $B$  carrying the same name is inaccessible in  $B$ .

( $\gamma$ ) When the algorithm run entered a block, a special internal structure called **block instance** was generated so that it carried entities local in the block. This structure was used only inside the compiled program and otherwise it was hidden. When the algorithm run left the block, the block instance was liquidated. But already according to the rules of ALGOL 60, it was possible that when the algorithm run was inside a certain block, it could enter the same block anew (it could happen e.g. in case of recursive calling procedures). In such a case each of the entry into the block caused the origin of a block instance and so more than one instance of the same block could exist contemporarily.

The whole apparatus of block handling was fully accepted into the simulation language SIMULA I (see above) and then into the object-oriented programming language SIMULA 67. For the last language, the following enlargements of the block apparatus appeared as natural and were introduced:

( $\delta$ ) As a local entity introduced for a block, a class can figure. Of course, the instances of such a class are meaningful only inside the block. The block with such a local class may be viewed as a “thinking” phase of the algorithm, using the class as a “private” concept.

( $\varepsilon$ ) As the life rules have form of algorithms, the blocks mentioned in ( $\delta$ ) were at disposal to represent thinking phases of the objects. More than one instance of the same class can exist and even be inside the same block (thinking phase). In such a case, for each of the class instances its “private” corresponding block instance is formed. Depending on the entities global for such a block, the class instances can differ and so different opinions of the “discussing” instances can be represented.

( $\zeta$ ) The blocks ordered in life rules similarly as in ( $\alpha$ ) can represent different thinking phases of the “life” of an instance. The homonymous classes local in such blocks may represent changes or especially development in the instance’s thinking

( $\eta$ ) The best description of a system  $S$  as a whole is a block  $B$  in that the classes of the elements of  $S$  figure as local classes. Let  $C1, C2, \dots$  be such local classes. In case the life rules of some of them, suppose  $C1$ , contains a block  $b$ , in which also classes like  $C1, C2, \dots$  are introduced as local,  $b$  might represent a “reflecting phase of  $C1$ ’s life”, i.e. as a phase when an instance of  $C1$  should something decide – often with respect to some future consequences – and it imagines (or – if it is a computer – simulates) the consequences of different variants of the decision and chooses the optimal one.

( $\theta$ ) In such a case,  $b$  is nested in  $B$  similarly as mentioned in ( $\beta$ ). Nevertheless, the so called dot notation, offered by SIMULA, enables surmounting the name conflicts and handling together the instances belonging to both the blocks (the main idea consists in transforming blocks to classes, some details are outlined in [17]). Note that this practice come to be fecund namely in case the studied systems are **anticipatory** ones in the weak sense [16, 18].

## Acknowledgement

This work has been supported by the Grant Agency of Czech Republic, grant reference no. 201/060612, name “Computer Simulation of Anticipatory Systems”.

## References

- [1.] MARKOV, A. A. (jun.): *Teorija algorifmov*. Academy of Sciences of USSR, Moscow, Leningrad, 1954.

- [2.] BACKUS, J. W. et al.: *Report on the Algorithmic Language ALGOL 60*. In Numerische Mathematik, Vol. 2, pp. 106-136, 1960.
- [3.] GORDON, G.: *A General Purpose Simulation Program*. In Proc. 1961 EJCC. MacMillan, New York, pp. 81-98, 1961.
- [4.] KNUTH, D. E., McNELEY, J. L.: *SOL – a symbolic language for general purpose systems simulation*. In IEEE Transactions on Electronic Computers, Vol. 13, pp. 401-409, 1964.
- [5.] KALINICHENKO, L. A.: *SLANG – Computer description and simulation-oriented experimental programming language*. In Simulation Programming Languages, North-Holland, Amsterdam, pp. 101-115, 1968.
- [6.] DAHL, O.-J., NYGAARD, K.: *SIMULA – A Language for Description of Discrete Event Systems: Introduction and User's Manual*. NCC Publ. No. 11, Norwegian Computing Center, Oslo, 1<sup>st</sup> edition 1965, 5<sup>th</sup> edition 1967
- [7.] DAHL, O.-J., NYGAARD, K.: *SIMULA – An ALGOL-Based Simulation Language*. In Communication of the ACM, Vol. 9, pp. 671-682, 1966.
- [8.] DAHL O.-J., NYGAARD, K.: *Class and Subclass Declarations*. NCC Publ. No. 93, Norwegian Computing Center, Oslo, 1967
- [9.] DAHL O.-J., MYHRHAUG, B., NYGAARD, K.: *Common Base Language*. NCC Publ. S-2. 1<sup>st</sup> edition 1968, 4<sup>th</sup> edition 1984.
- [10.] *SIMULA Standard as Defined by the SIMULA Standards Group 25<sup>th</sup> August 1986*. Simula a.s., Oslo, 1989
- [11.] DAHL, O.-J.: *Programming languages as tools for the formulations of concepts*. In Proc. 15<sup>th</sup> Scandinavian Congress, Oslo 1968. Lecture Notes in Mathematics no. 118, Springer, Berlin, pp. 18-29, 1970.
- [12.] MADSEN, O.-L., MØLLER-PEDERSEN, B., NYGAARD, K.: *Object-Oriented Programming in the Beta Programming Language*. Addison Wesley, Reading – Menlo Park – New York – Bonn – Amsterdam, 1993
- [13.] KLUSKA, V.: *Vytvoření programových prostředků pro diskrétní simulaci v jazyce JAVA* (Constructs of programming tools for discrete simulation in JAVA language). Master thesis in Czech, Ostrava University, Ostrava, 2007.
- [14.] HOLBAEK-HANSEN, E., HÅNDLYKKEN, P. and NYGAARD, K.: *System Description and the DELTA Language*. Norwegian Computing Center, Oslo, 1975.
- [15.] DAHL, O.-J.: *Relating a simulation model to an applicative specification*. In Modelling and Simulation – Proceedings ESM, Praha 1995. Society for Computer International, pp- 633-638, 1995.
- [16.] KINDLER, E.: *Object-Oriented Representations of Formal Theories as Tools for Simulation of Anticipatory Systems*. In Computing Anticipatory Systems CASYS 2005 – 7th International Conference, American Institute of Physics, Melville, New York, pp. 253-259, 2006.
- [17.] KINDLER, E., KŘIVÝ, I.: *A New SIMULA Class for Simulation*. In Proceedings of XXIX International Autumn Colloquium ASIS 2007 – Advanced Simulation of Systems [Acta MOSIS No. 110]. MARQ, Ostrava, 2007, pp. 7-12, 2007.
- [18.] KINDLER, E.: *Computer Models of Systems Containing Simulating Elements*. In Computing Anticipatory Systems CASYS 2000 - Fourth International Conference, Liege, Belgium [AIP Conference Proceedings Vol. 573]. American Institute of Physics, Melville, New York, pp. 390-399, 2001.

**Current address**

**Eugene Kindler, RNDr. PhDr. CSc, Professor**

Department of Informatics and Computers, Faculty of Sciences, University of Ostrava, 30. dubna no. 22, CZ-701 03 Ostrava, Czech Republic, phone 420 220 801 945,  
e-mail: [ekindler@centrum.cz](mailto:ekindler@centrum.cz)

**Ivan Křivý, RNDr. Ing. CSc, Professor**

Department of Informatics and Computers, Faculty of Sciences, University of Ostrava, 30. dubna no. 22, CZ-701 03 Ostrava, Czech Republic, phone 420 597 092 177,  
e-mail: [ivan.krivy@osu.cz](mailto:ivan.krivy@osu.cz)





# FRACTIONAL GENERALIZATION OF THE CLASSICAL VISCOELASTICITY MODELS

KISELA Tomas, (CZ)

**Abstract.** In this article we consider an application of the fractional calculus in the theory of viscoelasticity. First we give a brief survey of the important formulas of the fractional calculus and we mention an introduction into the classical viscoelasticity. Then we sketch main ideas of fractional viscoelasticity and finally we focus on a fractional model and derive its step response functions.

**Key words and phrases.** fractional calculus, viscoelasticity, Mittag-Leffler functions

*Mathematics Subject Classification.* Primary 26A33; Secondary 33E12.

## 1 Fundamentals of the Fractional Calculus

Fractional calculus is a mathematical discipline dealing with the so-called differintegrals (particularly fractional derivatives and fractional integrals). Differintegral is a natural generalization of the classical integral and derivative. There are many ways how to define it, but we are going to recall and employ here the Riemann-Liouville approach only. For more detailed information see [1] or [2].

**Definition:** Let  $a, T, \alpha$  be real constants ( $a < T$ ),  $n = \max(0, [\alpha] + 1)$  and let  $f(t)$  be an integrable function on  $[a, T)$ . For  $n > 0$  we additionally assume that  $f(t)$  is  $n$ -times differentiable on  $[a, T)$  except for a set of measure zero. Then the Riemann-Liouville differintegral of a function  $f(t)$  is defined for  $t \in \langle a, T)$  by the formula:

$$\mathbf{D}_a^\alpha f(t) = \frac{1}{\Gamma(n - \alpha)} \frac{d^n}{dt^n} \int_a^t (t - \tau)^{n - \alpha - 1} f(\tau) d\tau, \quad (1)$$

where  $\Gamma$  is the Euler Gamma function.

The main properties of differintegral we need in this article are their linearity and composition rules. The Riemann-Liouville definition uses only convolutions and classical derivatives, hence under some assumptions implied by the classical calculus we may write:

$$\mathbf{D}_a^\alpha \sum_{k=0}^{\infty} c_k f_k(t) = \sum_{k=0}^{\infty} c_k \mathbf{D}_a^\alpha f_k(t), \quad c_k \in \mathbb{R}. \quad (2)$$

Composition rules are a little more complicated and have the following form:

$$\mathbf{D}_a^\alpha (\mathbf{D}_a^\beta f(t)) = \begin{cases} \mathbf{D}_a^{\alpha+\beta} f(t), & \alpha \in \mathbb{R}, \beta \leq 0 \\ \mathbf{D}_a^{\alpha+\beta} f(t) - \sum_{k=1}^m \mathbf{D}_a^{\beta-k} f(t) \Big|_{t=a}^{\frac{(t-a)^{-\alpha-k}}{\Gamma(1-\alpha-k)}}, & \alpha \in \mathbb{R}, \beta \geq 0 \end{cases}. \quad (3)$$

Next let us look at the differintegrals of some functions. One of the most important functions used in fractional calculus is a power function because the differintegral of the power function remains still a power function. Hence, we may say that differintegration of the power function is analogous to multiplication with another power function. It holds

$$\mathbf{D}_a^\alpha (t-a)^\beta = \frac{\Gamma(\beta+1)}{\Gamma(\beta-\alpha+1)} (t-a)^{\beta-\alpha}, \quad \alpha \in \mathbb{R}, \beta > -1. \quad (4)$$

Other important function is the Mittag-Leffler function playing a very important role in the theory of linear fractional differential equations. It is defined by the relation

$$E_{\mu,\gamma}(t) = \sum_{k=0}^{\infty} \frac{t^k}{\Gamma(\mu k + \gamma)}, \quad \mu, \gamma \in \mathbb{R}, \mu > 0. \quad (5)$$

We can easily see that the case  $\mu = 1$  and  $\gamma = 1$  coincides with the classical exponential. The fractional analogy to the well-known equation  $y'(t) = y(t)$  and its solution  $y(t) = e^t$  is provided by the function  $y(t) = t^{\alpha-1} E_{\alpha,\alpha}(t^\alpha)$  solving the fractional differential equation  $\mathbf{D}_0^\alpha y(t) = y(t)$ . Hence we sometimes speak of generalized exponential instead of the Mittag-Leffler function.

More often we need a differintegral of a product of a function of Mittag-Leffler type and a power function, which can be calculated term by term via the formula (4)

$$\mathbf{D}_0^\alpha (t^{\beta-1} E_{\mu,\beta}(\lambda t^\mu)) = t^{\beta-\alpha-1} E_{\mu,\beta-\alpha}(\lambda t^\mu) \quad \alpha, \mu, \lambda \in \mathbb{R}, \beta > 0. \quad (6)$$

## The Laplace Transform

Similarly to the classical theory, the Laplace transform is a very powerful instrument for solving linear fractional differential equations with constant coefficients. Hence we need to know the Laplace image of the differintegral

$$\mathcal{L}\{\mathbf{D}_0^\alpha f(t), t, s\} = s^\alpha F(s) - \sum_{k=1}^n s^{n-k} \mathbf{D}_0^{\alpha-n+k-1} f(t) \Big|_{t=0}. \quad (7)$$

In the end of solving process we calculate the inverse transform of the solution. At this moment we appreciate a knowledge of the formula

$$\mathcal{L}\left\{t^{\alpha m + \beta - 1} E_{\alpha, \beta}^{(m)}(at^\alpha), t, s\right\} = \frac{m! s^{\alpha - \beta}}{(s^\alpha - a)^{m+1}}, \quad (8)$$

where  $E_{\alpha, \beta}^{(m)}(at^\alpha)$  denotes  $m^{th}$ -derivative of the Mittag-Leffler function according to the parameter  $a$ . This relation can be proved again term by term.

In this paper we will need also the formula for the Laplace transform of a power function

$$\mathcal{L}\{t^r\} = \frac{\Gamma(r+1)}{s^{r+1}}, \quad r > -1. \quad (9)$$

## 2 Classical Viscoelasticity

Viscoelasticity is a scientific discipline describing the material's behaviour via two physical quantities - stress  $\sigma(t)$  and strain  $\epsilon(t)$ . The stress is the average amount of force per unit area, the strain is the geometrical measure of deformation representing the relative displacement between particles in the material. Various models differ from each other in the way how to relate them.

To illustrate the behaviour of viscoelasticity models we consider a body in two situations. First we discuss the effect of sudden deformation - the strain function is described by the Heaviside unit step,  $\epsilon(t) = H(t)$ . The stress appropriate to this strain of the body is called the relaxation modulus and we denote it  $G(t)$ . The other interesting situation is the effect of sudden stress, i.e. the stress is described by the Heaviside unit step,  $\sigma(t) = H(t)$  and we are curious about the strain response called the stress creep compliance and denoted by  $J(t)$ . We usually use a common name "response functions" or "responses" for the functions  $G(t)$  and  $J(t)$ .

There are two basic models. The first one is the ideal solid, also called Hooke's element, with one parameter  $E$  - elastic constant. It is symbolized by a spring and it is described by the formula

$$\sigma(t) = E \epsilon(t). \quad (10)$$

The second one is the ideal fluid, called Newton's element, where the characteristic parameter is viscosity  $\eta$ . The symbol is a dashpot and the appropriate relation is

$$\sigma(t) = \eta \frac{d\epsilon(t)}{dt}. \quad (11)$$

For these models we get the relaxation moduli  $G_H(t)$ ,  $G_M(t)$  and stress creep compliances  $J_H(t)$ ,  $J_M(t)$  simply by substitution into the formulas (10) and (11) (see figures 1 and 2).

We can see that the relaxation modulus in Hooke's model is constant for  $t > 0$ , so there is no stress relaxation which is expected according to experiment. On the contrary, the stress relaxation occurs too fast (we may say immediately) in the case of Newton's model. Now let us look at the stress creep compliances. Hooke's model again keeps a constant value, Newton's one is more problematic because proportions of the body increase linearly to the infinity.

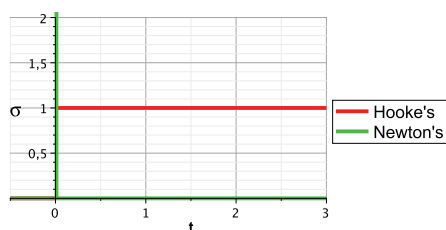


Figure 1: The relaxation moduli for ideal solid ( $E = 1$ ) and ideal fluid ( $\eta = 1$ ).

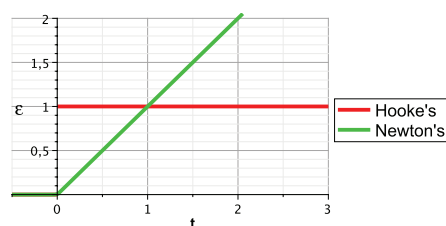


Figure 2: The stress creep compliances for ideal solid ( $E = 1$ ) and ideal fluid ( $\eta = 1$ ).

Obviously both models have serious problems with the description of the reality. In classical viscoelasticity we usually use serial (Maxwell's model) and parallel (Voigt's model) combinations of these two simple models for getting a better model. The responses of Maxwell's and Voigt's models are plot in figures 3 and 4. There is an improvement, but we are still unhappy about the relaxation modulus of Voigt's model and the stress creep compliance of Maxwell's model.

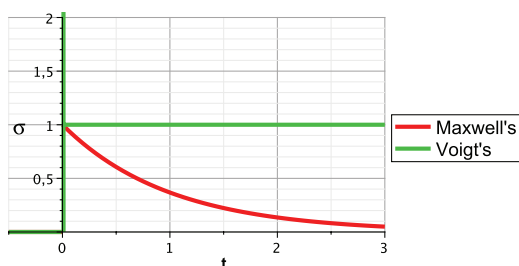


Figure 3: The relaxation moduli for Maxwell's and Voigt's models (constants  $E, \eta$  are equal to 1).

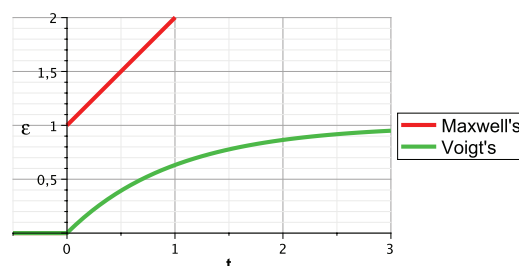


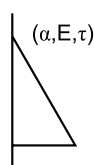
Figure 4: The stress creep compliances for Maxwell's and Voigt's models (constants  $E, \eta$  are equal to 1).

The quality of the model increases with the number of used elements. Nevertheless there are still problems like a finite value of the relaxation moduli at  $t = 0$  and the discontinuity of the stress creep compliance at  $t = 0$ . These difficulties can be reduced by adding more new elements.

### 3 Fractional Viscoelasticity Models

Fractional calculus brings new possibilities into modelling material's properties because the order of derivative plays a role of another parameter.

The idea is simply to take Hooke's and Newton's models and to realize that behaviour of real materials usually ranges between those two models. Hooke's model represents the zero derivative term and Newton's one corresponds with the first derivative term. Therefore a term with  $\alpha$ -derivative ( $0 \leq \alpha \leq 1$ ) is an intuitive generalisation and it is called Blair's model. There are three parameters - the order of differintegration  $\alpha$  and constants  $E, \tau$  forming one multiplicative constant (it is splitted just due to dependence on  $\alpha$ ). Its schematic symbol and its formula are



$$\sigma(t) = E\tau^\alpha \mathbf{D}_0^\alpha \epsilon(t), \quad \tau = \frac{\eta}{E}. \quad (12)$$

From the physical point of view we should choose a lower bound of differintegration in minus infinity because then we would include the whole history of the material. But we assume all quantities zero for  $t \leq 0$ , hence we may use the lower bound equal to zero.

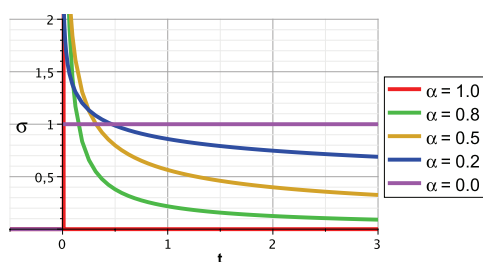


Figure 5: The relaxation moduli for Blair's model for various  $\alpha$  ( $E = 1$ ,  $\tau = 1$ ).

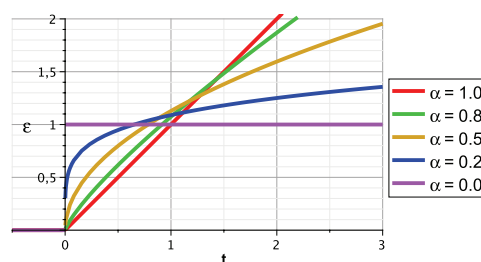


Figure 6: The stress creep compliances for Blair's model for various  $\alpha$  ( $E = 1$ ,  $\tau = 1$ ).

We can see from the graphs of the relaxation moduli and the stress creep compliances on figures 5 and 6 respectively that even this simple model provides the stress relaxation, the infinite value of the stress at  $t = 0$ , the continuity of the stress creep compliances and also their slower growing for greater  $t$ .

Of course, the description through presented power-laws is not always sufficient and we again use various combinations of Blair's elements. Now let us introduce one of them - two Blair's in series, so-called generalized Maxwell's model.

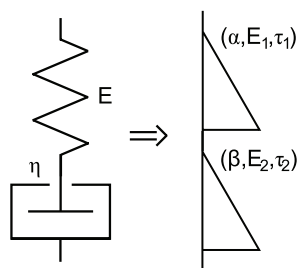


Figure 7: The schematical representation of the generalized Maxwell's model.

Obviously the stress is the same on both elements ( $\sigma_1(t) = \sigma_2(t) = \sigma(t)$ ) and the total strain is given by the sum  $\epsilon(t) = \epsilon_1(t) + \epsilon_2(t)$ . Hence we can write equations of this system in the form

$$\sigma(t) = E_1 \tau_1^\alpha \mathbf{D}_0^\alpha \epsilon_1(t), \quad (13)$$

$$\sigma(t) = E_2 \tau_2^\beta \mathbf{D}_0^\beta \epsilon_2(t). \quad (14)$$

Let us assume without loss of generality that  $\alpha > \beta$  and apply the operator  $\mathbf{D}_0^{\alpha-\beta}$  to the equation (14). According to the composition rule (3) we get

$$\mathbf{D}_0^{\alpha-\beta}\sigma(t) = E_2\tau_2^\beta \mathbf{D}_0^\alpha \epsilon_2(t) + E_2\tau_2^\beta \mathbf{D}_0^{\beta-1}\epsilon_2(t)\Big|_{t=0} \frac{t^{\beta-\alpha-1}}{\Gamma(\beta-\alpha)}. \quad (15)$$

The last term contains the initial condition for the strain but we postulated at the beginning of this section that all quantities are zero for  $t \leq 0$ . Hence this term disappears. The sum of the equations (13) and (15) gives the formula for this model

$$\frac{1}{E_1\tau_1^\alpha}\sigma(t) + \frac{1}{E_2\tau_2^\beta}\mathbf{D}_0^{\alpha-\beta}\sigma(t) = \mathbf{D}_0^\alpha \epsilon(t). \quad (16)$$

Let us note that the order connected with  $\sigma(t)$  is always less than the order incident to  $\epsilon(t)$  in the generalised Maxwell's model.

The derivation of the step responses  $G_{GM}(t)$ ,  $J_{GM}(t)$  uses the Laplace transform. First let us calculate the relaxation modulus  $G_{GM}(t)$ . We substitute  $\epsilon(t) = H(t)$  into (16) and according to formula (4) we arrive at

$$\frac{1}{E_1\tau_1^\alpha}\sigma(t) + \frac{1}{E_2\tau_2^\beta}\mathbf{D}_0^{\alpha-\beta}\sigma(t) = \frac{t^{-\alpha}}{\Gamma(1-\alpha)}.$$

Now we apply Laplace transform, particularly the formulas (7) and (9):

$$\frac{1}{E_1\tau_1^\alpha}\hat{\sigma}(s) + \frac{1}{E_2\tau_2^\beta}\left(s^{\alpha-\beta}\hat{\sigma}(s) - \mathbf{D}_0^{\alpha-\beta-1}\sigma(t)\Big|_{t=0}\right) = s^{\alpha-1}.$$

Again we can put  $\mathbf{D}_0^{\alpha-\beta-1}\sigma(t)\Big|_{t=0} = 0$  and then we express the Laplace image of the stress:

$$\hat{\sigma}(s) = \frac{E_2\tau_2^\beta s^{\alpha-1}}{s^{\alpha-\beta} + \frac{E_2\tau_2^\beta}{E_1\tau_1^\alpha}}.$$

The inverse Laplace transform of the  $\hat{\sigma}(s)$  is relaxation modulus  $G_{GM}(t)$  and we get it via relation (8):

$$G_{GM}(t) = E_2\tau_2^\beta t^{-\beta} E_{\alpha-\beta, 1-\beta} \left( -\frac{E_2\tau_2^\beta}{E_1\tau_1^\alpha} t^{\alpha-\beta} \right). \quad (17)$$

The situation is even more simple for derivation of stress creep compliance. We substitute  $\sigma(t) = H(t)$  into equation (16) and we can directly apply the differintegral  $\mathbf{D}_0^{-\alpha}$ . We use only formulas (3) and (4) to obtain

$$\begin{aligned} \frac{1}{E_1\tau_1^\alpha} + \frac{1}{E_2\tau_2^\beta} \frac{t^{-\alpha+\beta}}{\Gamma(1-\alpha+\beta)} &= \mathbf{D}_0^\alpha \epsilon(t), \\ \frac{1}{E_1\tau_1^\alpha} \frac{t^\alpha}{\Gamma(1+\alpha)} + \frac{1}{E_2\tau_2^\beta} \frac{t^\beta}{\Gamma(1+\beta)} &= \epsilon(t) - \mathbf{D}_0^{\alpha-1}\epsilon(t)\Big|_{t=0} \frac{t^{-\alpha-1}}{\Gamma(-\alpha)}. \end{aligned}$$

The last term again disappears due to zero initial condition, hence we obtain the stress creep compliance

$$J_{GM}(t) = \frac{1}{E_1 \tau_1^\alpha \Gamma(1 + \alpha)} t^\alpha + \frac{1}{E_2 \tau_2^\beta \Gamma(1 + \beta)} t^\beta. \quad (18)$$

Let us discuss possibilities provided by the generalized Maxwell's model. First let us keep  $\beta = 0$ , i.e. we examine the case of one Blair's and one Hooke's element in a series. It is evident from figure 8 that decreasing parameter  $\alpha$  provides a very effective way how to reduce the stress relaxation effect, and we also see that  $\beta = 0$  causes the boundedness of the relaxation moduli. In figure 9 we observe a flattening of the stress creep compliances with lowering parameter  $\alpha$ . The negative effect of zero  $\beta$  is discontinuity of those responses.

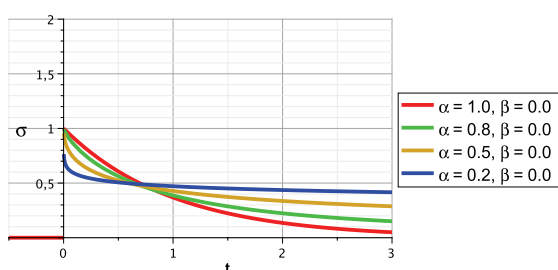


Figure 8: The relaxation moduli for the generalized Maxwell's models - influence of  $\alpha$  (all constants are equal to 1 except  $\alpha, \beta$ ).

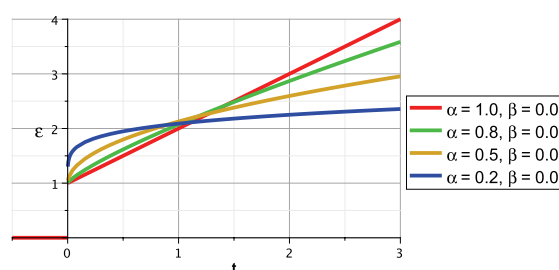


Figure 9: The stress creep compliances for the generalized Maxwell's models - influence of  $\alpha$  (all constants are equal to 1 except  $\alpha, \beta$ ).

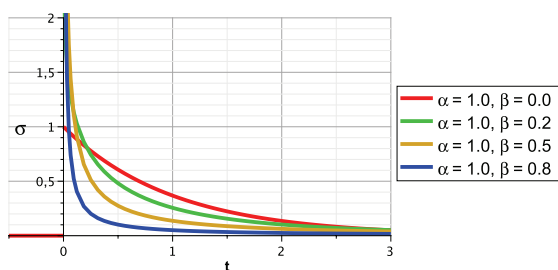


Figure 10: The relaxation moduli for the generalized Maxwell's models - influence of  $\beta$  (all constants are equal to 1 except  $\alpha, \beta$ ).

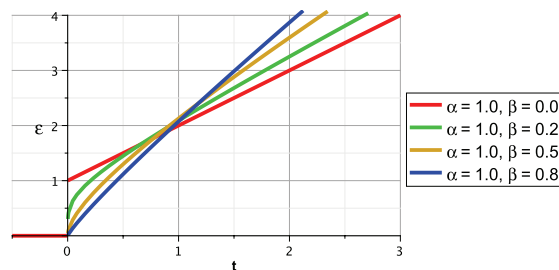


Figure 11: The stress creep compliances for the generalized Maxwell's models - influence of  $\beta$  (all constants are equal to 1 except  $\alpha, \beta$ ).

Now we fix the parameter  $\alpha$  on the value 1, so we are interested in the case of one Blair's and one Newton's element in series. The appropriate response functions are plotted in figures 10 and 11. Obviously increasing parameter  $\beta$  causes faster stress relaxation and its non-zero values bring unboundedness into a neighbourhood of the point  $t = 0$ . The non-zero parameter  $\beta$  also makes the stress creep compliances to be continuous. On the other hand the parameter  $\alpha = 1$  keeps the behaviour of those responses to be almost linear for  $t \gg 1$  which is not desired.

Clearly the limit values one, zero of the parameters  $\alpha, \beta$  respectively have a negative influence to response functions. Hence we will consider various values of  $\alpha, \beta$  but we will keep their difference constant (still  $\alpha > \beta$ ).

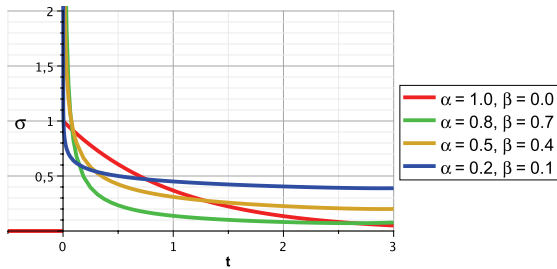


Figure 12: The relaxation moduli for the generalized Maxwell's models - influence of  $\alpha$  and  $\beta$  with constant difference (all constants are equal to 1 except  $\alpha, \beta$ ).

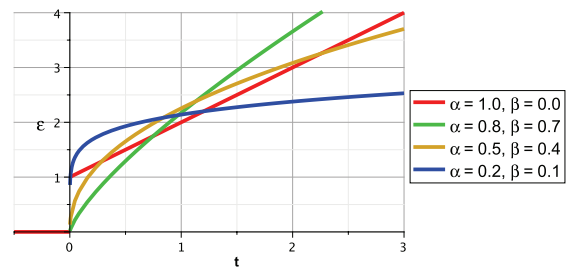


Figure 13: The stress creep compliances for the generalized Maxwell's models - influence of  $\alpha$  and  $\beta$  with constant difference (all constants are equal to 1 except  $\alpha, \beta$ ).

The appropriate curves are depicted in figures 12 and 13. The lines seem to be mixed up at the first glance, but there are some regularities. We see that for relaxation modulus decreasing  $\alpha$  causes slower stress relaxation for  $t \gg 1$ , we may say “fat tails”, whereas decreasing  $\beta$  effects faster stress fall for  $t \ll 1$ . The situation about the stress creep compliance is much more simple, because it behaves like  $\sim t^\alpha$  for  $t \gg 1$  like  $\sim t^\beta$  for  $t \ll 1$ .

Like in classical viscoelasticity, a greater number of parameters enables better adaptation of the shape of the curves to the reality. The difference is that we need a less number of elements with fractional models for qualitatively suitable results.

## 4 Conclusions

In this paper we discussed one of the most important application of the fractional calculus, the theory of viscoelasticity. In general the fractional calculus provides a very interesting instrument for modelling because the order of differintegration plays a role of a new parameter. That is the “macroscopic” reason why we use this theory in viscoelasticity. On the other side there exists a microscopic point of view which also arrives to fractional models as we can see e.g. in [3].

We introduced the complete macroscopic derivation of the generalized Maxwell's model, and we described its behaviour, particularly the effects of the change of parameters  $\alpha, \beta$  on the response functions  $G(t)$  and  $J(t)$ .

In [3] there is also mentioned that more general fractional models are very successful in describing viscoelastic properties of many polymeric materials, especially their relaxation behaviour. Some information about this phenomena and also about another application of the fractional calculus can be found e.g. in [1], [2].

## Acknowledgement

The paper was supported by the research plan MSM 0021630518 “Simulation modelling of mechatronic systems” of the Ministry of Education, Youth and Sports of the Czech Republic.



## References

- [1] PODLUBNY, Igor. *Fractional Differential Equations*. United States: Academic Press c1999. 340 p. ISBN 0-12-558840-2.
- [2] KILBAS, Anatoly A., SRIVASTAVA, Hari M., TRUJILLO, Juan J. *Theory and Applications of Fractional Differential Equations*. John van Mill. Netherlands: Elsevier, 2006. 523 p. ISBN 978-0-444-51832-3.
- [3] HILFER, R. *Applications of Fractional Calculus in Physics*. Singapore: World Scientific 2000. 463 p.

## Current address

**Tomas Kisela, Ing.**

Faculty of Mechanical Engineering

Brno University of Technology

Technicka 2896/2

616 69 Brno

e-mail: ykisel00@stud.fme.vutbr.cz



## STABILISATION OF MEAN AND VARIANCE FOR NONSTATIONARY PROCESSES

**KVAPIL David, (CZ)**

**Abstract.** Nonstationary processes occur in stochastic analysis of technological data. This problem can be solved by several methods. We resume some approaches to the task of nonstationary process and shortly illustrate process of stochastic analysis in technometrics. We will construct GARCH model for technological data and demonstrate its creation in MATLAB.

**Keywords.** Stochastic modeling, Box – Jenkins (S)AR(I)MA model, nonstationary process, GARCH model

### 1 Introduction

Analysis of technological data indicates a high-frequency time series with changeable variance, we can speak about volatility. There are many problems in practice classic linear (S)AR(I)MA models (Box – Jenkins methodology), which allow only correlation dependence. The variability relates to the autocorrelation. We can understand the change of volatility as the change of the time series regime which is determined by different deterministic and unsystematic factors.

For empirical daily time series there are usually not satisfied conditions of linear modelling (homoskedasticity and normality). After the graphical analysis we can see that the data comes from leptokurtic distribution (fat tails, excess kurtosis). The first conception was proposed by Engle (1984) – his ARCH model supposed conditional variance. The requirement of normality was preserved.

### 2 Nonstationary mean and variance

For nonstationary processes in mean we distinguish a deterministic trend and a stochastic trend. For the deterministic trend nonstationariness is perceived as a function of the time. For modeling we use a polynomial or periodic trend, respectively

$$f(t) = \beta_0 + \beta_1 t + \dots + \beta_d t^d, \quad f(t) = \mu + \sum_{j=1}^p (\alpha_j \cos \lambda_j t + \beta_j \sin \lambda_j t).$$

For the Box – Jenkins ARMA models

$$ARMA(p, q): Y_t - \varphi_1 Y_{t-1} - \dots - \varphi_p Y_{t-p} = \varepsilon_t + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q},$$

where  $\varepsilon_t \sim WN(0, \sigma_\varepsilon^2)$ , expressed using backshift operator  $\Phi(B)Y_t = \Theta(B)\varepsilon_t$ , is requested causal process, i.e. each of the roots of polynomial  $\Phi(z) = 1 - \varphi_1 z - \dots - \varphi_p z^p$  is outside the unit circle. If any root is situated on the unit circle, we speak about nonstationary process with the stochastic trend. If any root is situated inside the unit circle, we speak about the nonstationariness of explosive type. [9]

We can reduce the stochastic trend using the difference operator

$$\begin{aligned} \Delta Y_t &= Y_t - Y_{t-1} = (1 - B)Y_t, \\ \Delta^2 Y_t &= \Delta(\Delta Y_t) = \Delta(Y_t - Y_{t-1}) = Y_t - 2Y_{t-1} + Y_{t-2} = (1 - B)^2 Y_t, \\ \Delta^d Y_t &= Y_t - \binom{d}{1} Y_{t-1} + \binom{d}{2} Y_{t-2} - \dots + (-1)^d Y_{t-d} = (1 - B)^d Y_t. \end{aligned}$$

Nonstationary process with the stochastic trend is called integrated ARIMA( $p, d, q$ ) model, we can write

$$ARIMA(p, d, q): \Phi(B)(1 - B)^d Y_t = \Theta(B)\varepsilon_t.$$

Note that the number  $d$  may not be integer; then  $d$  is called the fractional parameter, we work with the fractional difference and we have the fractional integrated process ARFIMA( $p, d, q$ ). [3]

The non-stable variance process can be reduced by the Box – Cox transformation or power transformation, respectively (for  $Y_t > 0$ )

$$Z_t = \begin{cases} (Y_t^\lambda - 1)/\lambda & \text{for } \lambda \neq 0, \\ \ln(Y_t) & \text{for } \lambda = 0, \end{cases} \quad Z_t = \begin{cases} Y_t^\lambda & \text{for } \lambda \neq 0, \\ \ln(Y_t) & \text{for } \lambda = 0. \end{cases}$$

If random variables  $Y_t$  are not positive, we can use the following transformations

$$Z_t = \begin{cases} \frac{(Y_t + a)^\lambda - 1}{\lambda} & \text{for } \lambda \neq 0, \\ \ln(Y_t + a) & \text{for } \lambda = 0, \end{cases} \quad Z_t = \begin{cases} \operatorname{sgn}(Y_t) \frac{|Y_t|^\lambda - 1}{\lambda} & \text{for } \lambda \neq 0, \\ \operatorname{sgn}(Y_t) \ln(Y_t) & \text{for } \lambda = 0. \end{cases}$$

Note that for  $Y_t < 0$  (or near zero) we can make any transformation only with the knowledge in significant risk for degradation the time series and incredible final model. Beyond the power transformation is not continuous for  $\lambda \rightarrow 0$ ; it is necessary to keep away from small nonzero  $\lambda$ . [8] The estimation of the transformation parameter is performed using the maximum of logarithm of likelihood function

$$l^*(\lambda) = -\frac{n}{2} \ln(\hat{\sigma}^2(\lambda)) + (\lambda - 1) \sum_{i=1}^n \ln y_i.$$

All  $\lambda$  satisfy

$$l^*(\lambda) \geq D_\alpha = l^*(\hat{\lambda}) - \frac{1}{2} \chi_{1-\alpha}^2(1)$$

are situated in the confidence interval. [9]

### 3 GARCH modelling

An alternative approach is the modelling of processes with changeable regime – TAR (Threshold Autoregressive), MSW (Markov Switching) or ARCH/GARCH (Generalized Autoregressive Conditional Heteroskedasticity) models.

The GARCH volatility model employs conditional mean and conditional variance. The conditional variance is a linear function of values  $\varepsilon_{t-1}^2, \varepsilon_{t-2}^2, \dots, \varepsilon_{t-p}^2$  for linear volatility models. Nonlinear models are able to represent certain asymmetric events (e.g. leverage effects). Basic linear models are ARCH( $q$ ) and GARCH( $p, q$ ). ARCH( $q$ ) model is

$$\sigma_t^2 = \omega + \alpha_1 \varepsilon_{t-1}^2 + \alpha_2 \varepsilon_{t-2}^2 + \dots + \alpha_q \varepsilon_{t-q}^2,$$

where  $\omega > 0$  and  $\alpha_1, \dots, \alpha_q \geq 0$ . GARCH( $p, q$ ) model is

$$\sigma_t^2 = \omega + \alpha_1 \varepsilon_{t-1}^2 + \alpha_2 \varepsilon_{t-2}^2 + \dots + \alpha_q \varepsilon_{t-q}^2 + \beta_1 \sigma_{t-1}^2 + \beta_2 \sigma_{t-2}^2 + \dots + \beta_p \sigma_{t-p}^2,$$

where  $p \geq 0, q > 0, \omega > 0, \alpha_i \geq 0$  for  $i = 1, \dots, q, \beta_j \geq 0$  for  $j = 1, \dots, p$ .

Other linear models are IGARCH, FIGARCH, GARCH-M, nonlinear models are EGARCH, IEGARCH, FIEGARCH, GJR-GARCH, QGARCH, SV model, etc. [3] [4] [6]

The general ARMAX( $p, q, n$ ) model for the conditional mean

$$Y_t = C + \sum_{i=1}^p \varphi_i Y_{t-i} + \varepsilon_t + \sum_{j=1}^q \theta_j \varepsilon_{t-j} + \sum_{k=1}^n \beta_k A_{t,k},$$

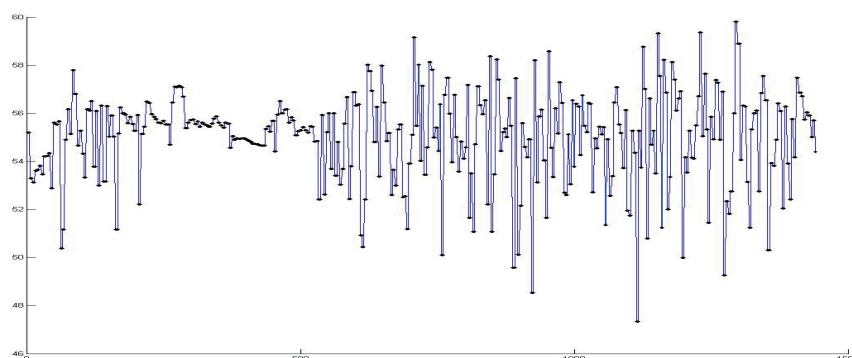
where  $A_{t,k}$  is an explanatory regression matrix in which each column is a time series, applies to all variance models.

The process of construction of volatility models is as follows: [3]

- i) Fitting linear or nonlinear level model is created for the time series.
- ii) Null hypothesis of conditional homoskedasticity is tested against alternative hypothesis of conditional heteroskedasticity of linear or nonlinear type.
- iii) Parameters of linear or nonlinear selected model of conditional heteroskedasticity are estimated.
- iv) Fitness of selected model is verified by diagnostic tests.
- v) The model is modified if it is necessary.
- vi) The model is used for description or prediction.

#### 4 Analysis of technological data

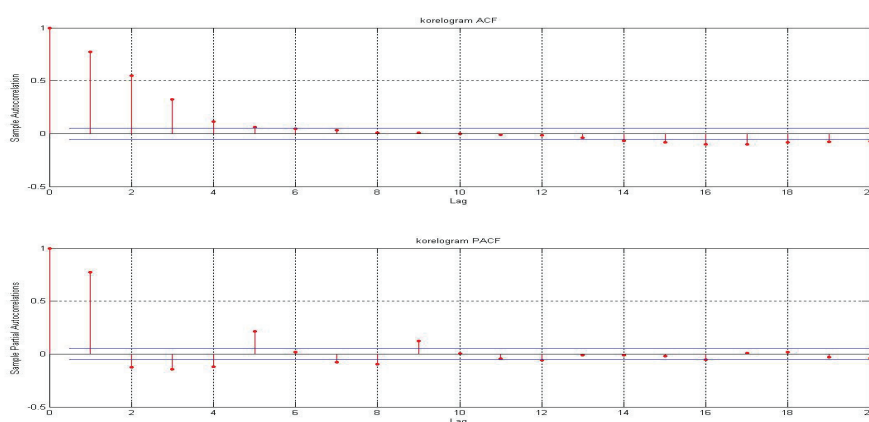
We shortly demonstrate the stochastic analysis of technological data from the thermal power station. In the next figure there is the process of temperature of output warm servis water during one day.



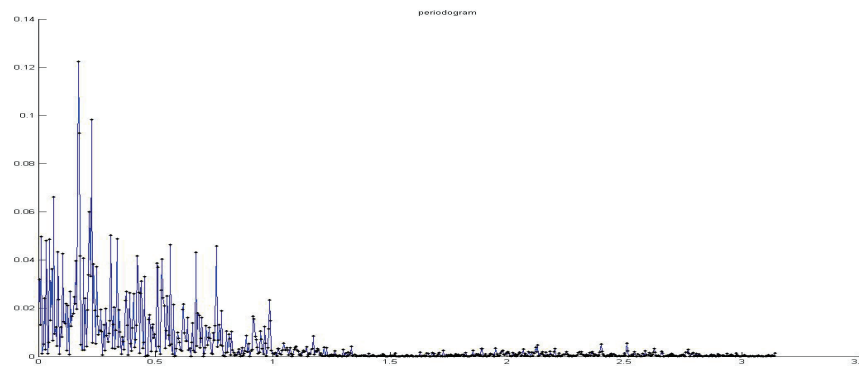
For the data vector (denoted A) we create MATLAB script BJmodel.m. It realizes the identification of Box – Jenkins model. We obtain

```
>> BJmodel(A)
vyberovy prumer: 55.0193
vyberovy rozptyl: 3.60564
hodnoty FPE kriteria pro jednotlive ARMA modely:
1.0e+003 *
    3.030731706022    0.909115365313    0.324244052345    0.166364044939
    0.001631776099    0.001634188044    0.001636436964    0.001638666986
    0.001634048368    0.001636455241    0.001638704943    0.001637861626
    0.001637473969    0.001639879685    0.001641482159    0.001642407768
```

From generated matrix of values of FPE criterion we can choose ARMA(1,0) because the value of the element in position (2,1) is minimal. The correlogram is in the next figure.



The periodogram is in the next figure, there is not indicated any significant frequency.

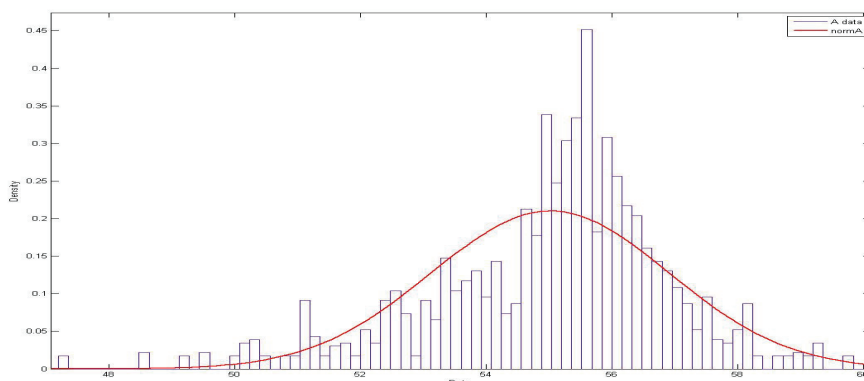


We create MATLAB script ARMAModel.m for verification ARMA(1,0) model. We obtain

```
>> ARMAModel(A,1,0)
Discrete-time IDPOLY model:  $A(q)y(t) = e(t)$ 
 $A(q) = 1 - 0.9731 (+0.01949) q^{-1}$ 
Estimated using ARMAX from data set yc
Loss function 1.58739 and FPE 1.5896
Sampling interval: 1
```

```
hodnota Portmonteau statistiky:  $Q = 251.591$ 
kriticka hodnota : krit = 52.1923
```

The value of Portmonteau statistic is  $Q = 251.591$  and it is much higher then the critical value of the test  $k = 52.1923$ . We create the histogram of measured data and compare it with the normal distribution.



We can see that we have data from the leptokurtic distribution. Now we perform Jarque – Bera test (jbtest) of normality in MATLAB.

```
>> [h,p,jbstat,krit] = jbtest(A, 0.05);
>> [h,p,jbstat,krit]
ans =
    1.00000000000000    0.00100000000000    227.0430989943586    5.9495166571438
```

The test rejects null hypothesis at the 5% significance level, our data does not come from the normal distribution.

We create MATLAB script testnezprvku.m (test of independence). It performs the test of autocorrelation of sample elements.

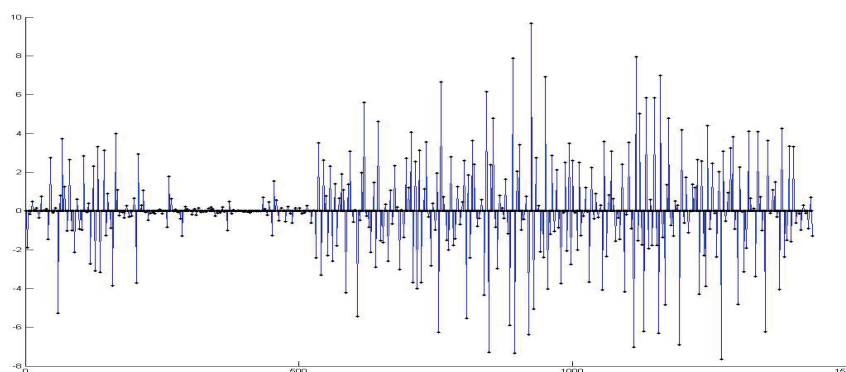
```
>> testnezprvku(A)
*** H0: neni autokorelace ***
H0 zamitame - existuje autokorelace
testovaci statistika: 61.8415
hranice kritickeho oboru:
-1.961611612000120 1.961611612000120
```

The test rejects null hypothesis at the 5% significance level, the autocorrelation exists in our sample.

Now we differentiate our sample and we will work with the series of the first differences.

```
>> A1=diff(A);
```

On the next figure (plot of A1) we can see the clusters of volatility.



We use archtest in MATLAB for test of presence ARCH effects (heteroscedasticity, leptokurtic distribution, leverage effect).

```
>> [H,p,stat,krit] = archtest(A1-mean(A1),[1 2 3 4 5 10 15 20]', 0.05);
>> [H,p,stat,krit]
```

ans =

1.000000000000	0.025741377174	4.9733075496981	3.841458820694
1.000000000000	0.005051016933	10.576331366106	5.991464547108
1.000000000000	0.010148249556	11.313022268465	7.814727903251
1.000000000000	0	139.05570911226	9.487729036781
1.000000000000	0	152.31277873385	11.07049769351
1.000000000000	0	155.15336092862	18.30703805327
1.000000000000	0	159.15625510344	24.99579013972
1.000000000000	0	161.92284873141	31.41043284423

The test rejects null hypothesis at the 5% significance level, the conditional variance exists in our model for the values  $q = 1, 2, 3, 4, 5, 10, 15$  and 20.

Now we use likelihood ratio hypothesis test for estimated GARCH( $p, q$ ) models.

```
>> spec11 = garchset('P',1,'Q',1,'Display','off');
>> spec21 = garchset('P',2,'Q',1,'Display','off');
```



```
>> spec12 = garchset('P',1,'Q',2,'Display','off');
>> spec22 = garchset('P',2,'Q',2,'Display','off');
>> [coeff11,errors11,LLF11] = garchfit(spec11,A1);
>> [coeff12,errors12,LLF12] = garchfit(spec12,A1);
>> [coeff21,errors21,LLF21] = garchfit(spec21,A1);
>> [coeff22,errors22,LLF22] = garchfit(spec22,A1);

>> [H,pValue,Stat,CriticalValue] = lratiotest(LLF12,LLF11,1,0.05);
>> [H,pValue,Stat,CriticalValue]
ans =
    1.000000000000000    0.000000000000161  49.9072579043968  3.8414588206941

>> [H,pValue,Stat,CriticalValue] = lratiotest(LLF12,LLF21,1,0.05);
>> [H,pValue,Stat,CriticalValue]
ans =
    1.000000000000000    0.000000000000161  49.9072577831202  3.8414588206941
```

The more formal approach to the choice of order of the model is using certain criterion function. The most common and widely used are FPE (Final Prediction Error, 1969), AIC (Akaike Information Criterion, 1974), AICC (corrected AIC), BIC (Bayesian Information Criterion, 1978), SBC (1978) etc. [9]

We create MATLAB script GARCHmodel.m which generates matrices of values of AIC and BIC criteria, respectively.

```
>> GARCHmodel(A1,2)
hodnoty AIC kriteria pro jednotlivé GARCH modely:
    1.0e+003 *
    4.220753025055129    4.172845767150732
    4.222753024933852    4.172967839307103
hodnoty BIC kriteria pro jednotlivé GARCH modely:
    1.0e+003 *
    4.241839839882679    4.199204285685170
    4.249111543468290    4.204598061548427
```

Now we have GARCH(1,2) model as the best option. We prepare MATLAB structure for this model and then we estimate all parameters of the model.

```
>> spec = garchset('P', 1, 'Q', 2)
>> [coeff,errors,LLF,eFit,sFit] = garchfit(spec,A1);
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
Diagnostic Information

Number of variables: 5

Functions
Objective:          garchllfn
Gradient:           finite-differencing
Hessian:            finite-differencing (or Quasi-Newton)
Nonlinear constraints: armanlc
Gradient of nonlinear constraints: finite-differencing

Constraints
```

```

Number of nonlinear inequality constraints:    0
Number of nonlinear equality constraints:    0
Number of linear inequality constraints:     1
Number of linear equality constraints:       0
Number of lower bound constraints:          5
Number of upper bound constraints:          5

```

```

Algorithm selected
  medium-scale

```

```

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
End diagnostic information

```

```

>> garchdisp(coeff,errors)
  Mean: ARMAX(0,0,0); Variance: GARCH(1,2)
  Conditional Probability Distribution: Gaussian
  Number of Model Parameters Estimated: 5

```

Parameter	Value	Standard Error	T Statistic
-----	-----	-----	-----
C	-0.011966	0.01614	-0.7413
K	0.0015184	8.0975e-005	18.7519
GARCH(1)	0.96139	0.0012091	795.1219
ARCH(1)	0	0.010684	0.0000
ARCH(2)	0.03861	0.01047	3.6876

Finally, we have GARCH(1,2) model

$$Y_t = -0.011966 + \varepsilon_t,$$

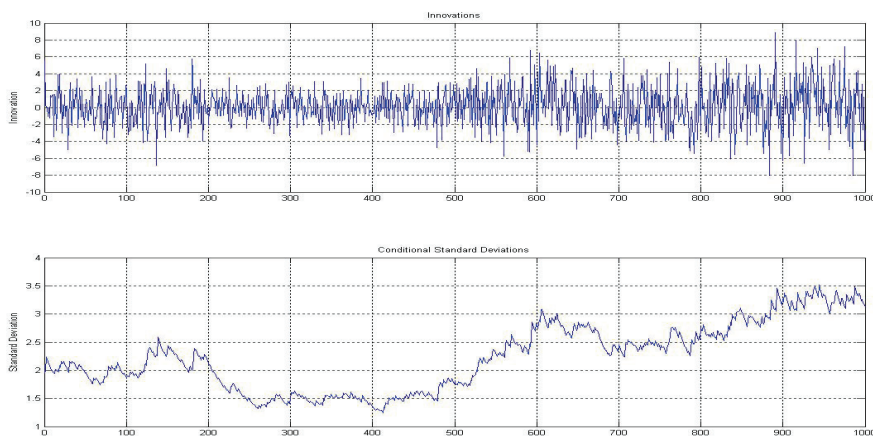
$$\sigma_t^2 = 0.0015184 + 0.96139\sigma_{t-1}^2 + 0.03861\varepsilon_{t-2}^2.$$

Now we simulate the data by the model above. The residua and conditional standard deviation are plotted in the next figure.

```

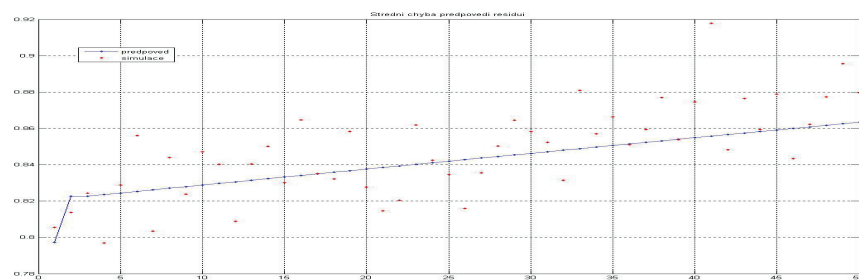
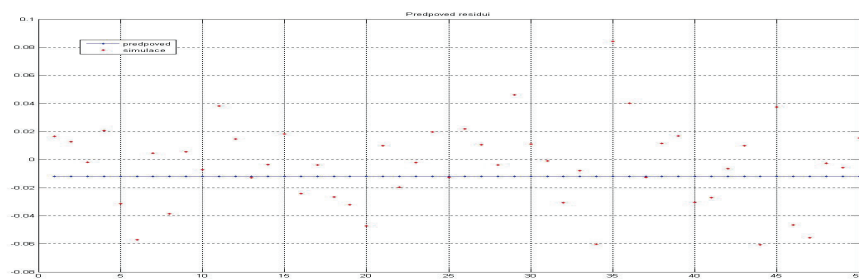
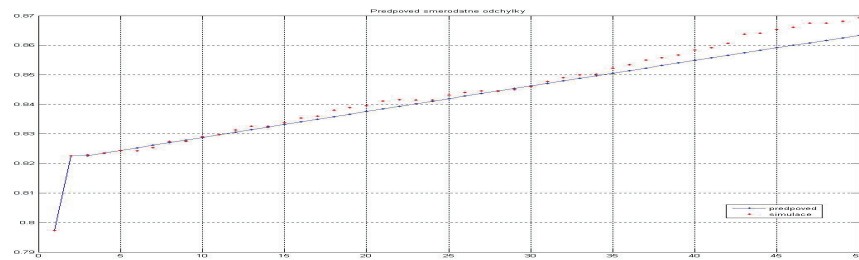
>> [e,s] = garchsim(coeff,1000);
>> garchplot(e,s)

```



Lastly we will predict with the horizon 50 samples and compare it with the results of simulation of our model. We will use Monte Carlo simulation of forecast.

```
>> hor=50;  
>> [sigmaForecast,meanForecast,sigmaTotal,meanRMSE] =  
garchpred(coeff,A1,hor);  
>> nPaths=1000;  
>> [eSim,sSim,ySim] =  
garchsim(coeff,hor,nPaths,0,[],[],eFit,sFit,A1(end));
```



## 5 Conclusion

Our result is not a product of a completed project; it is only a part of running partial problem of modeling and prediction of consumption of heat energy. GARCH models are commonly used and often applied in econometrics. Our aim is to demonstrate how to use these stochastic models in technometrics.

## References

- [1.] ANDĚL, J.: *Matematická statistika*, SNTL/ALFA, Praha, 1978
- [2.] ANDĚL, J.: *Statistická analýza časových řad*, SNTL, Praha, 1976
- [3.] ARLT, J. – ARLTOVÁ, M.: *Finanční časové řady*, Grada, Praha, 2003
- [4.] ARLT, J. – ARLTOVÁ, M.: *Konstrukce předpovědi na základě modelů GARCH*, Acta oeconomica pragensia 10: (7), VŠE, Praha, 2002
- [5.] ARLT, J.: *Moderní metody modelování ekonomických časových řad*, Grada, Praha, 1999
- [6.] ARLT, J. – Radkovský, Š.: *Význam modelování a předpovídání volatility časových řad pro řízení ekonomických procesů*, Politická ekonomie 48: (1), VŠE, Praha, 2000
- [7.] BUDÍKOVÁ, M. – LERCH, T. – MIKOLÁŠ, Š.: *Základní statistické metody*, skriptum MU, Brno, 2005
- [8.] CIPRA, T.: *Analýza časových řad s aplikacemi v ekonomii*, SNTL, Praha, 1986
- [9.] FORBELSKÁ, M.: *Stochastické modelování jednorozměrných časových řad I, II*, skriptum MU, Brno, 2007
- [10.] HEBÁK, P. a kol.: *Vícerozměrné statistické metody 1, 2, 3*, Informatorium, Praha, 2005, 2007
- [11.] LIKEŠ, J. – MACHEK, J.: *Matematická statistika*, MVŠT – sešit IV, SNTL, Praha, 1983
- [12.] LIKEŠ, J. – MACHEK, J.: *Počet pravděpodobnosti*, MVŠT – sešit X, SNTL, Praha, 1981
- [13.] MELOUN, M. – MILITKÝ, J.: *Statistická analýza experimentálních dat*, Academia, Praha, 2004
- [14.] RAO, R. C.: *Lineární metody statistické indukce a jejich aplikace*, Academia, Praha, 1978

## Current address

**David KVAPIL, RNDr.**

UNIS, a.s.

Department of Optimization and Advanced Control

Jundrovská 33, 624 00 Brno, CZ

+420 541 515 390

e-mail:dkvapil@unis.cz

## PARALLEL POSIX THREADS BASED ANT COLONY OPTIMIZATION USING ASYNCHRONOUS COMMUNICATIONS

LUCKA Maria, (SK), PIECKA Stanislav, (SK)

**Abstract.** In this paper we study parallel Posix threads based implementation of Ant Colony Optimization for solving the Vehicle Routing Problem. The algorithm is based on a homogeneous multi-colony approach and uses asynchronous communication in finding solutions. The aim of such approach is to examine potential advantage of executing parallel algorithm within the multicore processor environment. We analyze the effect of proposed method on the quality of solution with respect to execution and communication time.

**Key words and phrases.** Ant Colony Optimization, Parallel Metaheuristic, Vehicle Routing Problem, POSIX threads.

*Mathematics Subject Classification.* Parallel computation 65Y05; Transportation 90B06; Parallel algorithms 68W10.

### 1 Introduction

The Vehicle Routing Problem (VRP) is a combinatorial problem with numerous applications in telecommunication, transportation, logistics, etcetera. The majority of these applications belong to NP-hard problems, where to find the optimal solution requires in the worst case the exponential time. The VRP problem incorporate the construction of a set of vehicle tours that start and end at a depot and satisfy the demands of a set of customers. Each customer is served exactly once and both vehicle capacities and maximum tour lengths cannot be violated. For these problems no exact polynomial solutions are available and all known solutions so far are obtained by heuristic or metaheuristic algorithms. Especially metaheuristic methods seem to produce quality solutions in shorter calculation times.

Very successful method for solving the VRP problem is the Ant Colony Optimization (ACO) method, developed by [9]. It was inspired by behavior of real ants that deposit pheromone on ground to inform other ants about paths which should be followed by them. The idea of the algorithm is taken from the observation of the ant behavior. Each ant in the colony builds repeatedly its own solutions that are dependent on a given instance, a joint memory and an heuristic information. After all ants in a colony have found their own solutions of given combinatorial problem, the best solutions are selected and used for updating of the common memory. In computer implementation, the ACO method is represented by repeatedly called procedures which create solution by exploring fully connected graph of customers. After the solution is build, the pheromone matrix is updated according to achieved quality of solution.

There are more applications of ACO on VRP variants published by scientists. Applications based on basic VRP (CVRP respectively) can be found in [2], [13], [7]. ACO on VRP with Time Windows (VRPTW) are described in [10] and VRP with Pick-ups and Deliveries in [8]. There are more applications of ACO on VRP extensions listed in [11], where detailed analysis of parallel ACO methods used for solving the VRP problem can be found, however all of them are NP-Hard.

Many scientists have aimed to find parallel variants of the ACO for several reasons. The first is the fact that searching for new solution makes the work of each ant in a colony independent from other ants. The next reason is the fact that the time needed for solving the optimization problem is proportional to the size of instance depending so on the number of customers. Even finding of an optimal solution is for larger instances impossible due to the NP-hardness, the time needed for solving the real problems increases exponentially.

In this study we have chosen to use Savings based ACO algorithm for VRP as described in [7]. We study the VRP problem for large instances [4], [12] of basic VRP also called Capacitated Vehicle Routing Problem. Instances of such scale still don't have exactly calculated optimal solutions due to the problem's NP-Hardness. Our aim was to design a parallel algorithm for ACO suitable for multicore architectures. We have modified the parallel algorithm used in [7] for the multi-colony approach. We have used coarse-grained parallelization strategy, whereby each computing thread has assigned exactly one colony of ants. We suppose, that all colonies have the same behavior and are homogeneous. They cooperate in finding the best solution in an asynchronous way. The threads are divided into groups according to the number of cores in a computational node of the computer. They work in parallel and exchange the best known solutions so far within a group residing in a node. For communication within a group common shared memory is used. The groups of threads belonging to different nodes exchange the solutions via shared files across the network. Those groups are based on cluster topology and one group of threads corresponds to one node of the cluster. The communication between groups of threads residing on different nodes is performed over the network by means of usage of shared files.

The paper is organized as follows. Next section brings the formulation the VRP problem. In Section 3 we describe shortly parallelization strategies for ACO, and outline the algorithm used for implementation with Posix threads. Section 4 brings gained computational results and shows dependence of the solution quality on the number of threads, whereby the execution and communication time are also presented. We conclude with several remarks and outlooks concerning the future work.

## 2 Problem formulation of the Vehicle Routing Problem

The Vehicle Routing Problem (VRP) can be according to [4] described as follows: Let  $G = (V, E, c)$  be a complete graph, with  $n + 1$  nodes  $(v_0, \dots, v_N)$  corresponding to the customers  $i = 1, \dots, N$  and the depot  $i = 0$ , and the edge set  $((v_i, v_j) \in E \forall v_i, v_j \in V)$ . With each edge  $(v_i, v_j) \in E$  is associated a non-negative weight  $c_{ij}$ , which refers to the travel costs between nodes  $v_i$  and  $v_j$  and a non-negative weight  $t_{ij}$ , which refers to the distance between the nodes. Furthermore, with each node  $v_i, i = 1, \dots, N$  is associated a non-negative demand  $d_i$ , which has to be satisfied, as well as a service time  $\delta_i$ . The service time at the depot is set to  $\delta_0 = 0$ . At the depot a fleet of size  $K$  is available, where each vehicle has a capacity of  $Q^k$  and the maximum driving time for each vehicle is  $T^k$ .

Let  $x_{ij}^k$  denote the binary decision variables with the following interpretation:

$$x_{ij}^k = \begin{cases} 1 & \text{if vehicle } k \text{ visits node } v_j \\ & \text{immediately after node } v_i \\ 0 & \text{otherwise.} \end{cases}$$

Then the objective can be written as

$$\text{minimize} \sum_{i=0}^N \sum_{j=0}^N \sum_{k=1}^K c_{ij} x_{ij}^k \quad (1)$$

under the following restrictions

$$\sum_{i=1}^N \sum_{j=1}^N x_{ij}^k d_i \leq Q^k \quad 1 \leq k \leq K \quad (2)$$

$$\sum_{i=0}^N \sum_{j=0}^N x_{ij}^k (t_{ij} + \delta_i) \leq T^k \quad 1 \leq k \leq K \quad (3)$$

$$\sum_{i=0}^N x_{ij}^k - \sum_{l=0}^N x_{jl}^k = 0 \quad 1 \leq k \leq K, 0 \leq j \leq N \quad (4)$$

$$\sum_{i=0}^N \sum_{k=1}^K x_{ij}^k = \begin{cases} 1 & 1 \leq j \leq N \\ K & j = 0 \end{cases} \quad (5)$$

$$\sum_{i \in S} \sum_{j \in S} x_{ij}^k \leq |S| - 1 \quad \forall S \subseteq \{1, \dots, N\}, 1 \leq k \leq K \quad (6)$$

$$x_{ij}^k \in \{0, 1\} \quad 1 \leq k \leq K, 0 \leq i, j \leq N \quad (7)$$

The objective (1) is to minimize the total travel costs. Constraints (2) ensure that no vehicle is overloaded. Constraints (3) require that the maximum driving time for each vehicle is respected. Constraints (4) ensure that if a vehicle visits a customer it also leaves the customer. Constraints

(5) require that all customers are visited once, and that the depot is left  $K$  times. Sub-tour elimination is ensured through constraints (6). Finally, constraints (7) are the usual binary constraints.

### 3 Parallelization of ACO

The application of distributed computing on ACO means that several processes or threads work simultaneously on several processors to obtain common solution with the best solution quality. This can be achieved by functional or domain decomposition of problem and distribution over processors. Functional parallelism is characterized by several tasks working over same data. Domain decomposition is characterized by splitting the data into several smaller parts which are calculated separately. Tasks are typically cooperating when calculating solution, therefore some communication is required. In case that tasks are synchronizing themselves or not we can distinguish between synchronous and asynchronous communication model. For more details concerning the parallelization of ACO for VRP see [7], and [6]. In general, there are three possibilities of parallelization of ACO: fine-grained parallelization, coarse-grained parallelization and mixed parallelization. The fine-grained parallelization is a low-level parallelization achieved by splitting a colony into several sub-colonies that are processed in parallel. By coarse-grained parallelization the parallel search is computed by using several homogeneous colonies. The mixed parallelization is a combination of fine and coarse-grained parallelization.

In our work we have rewritten parallel algorithm using MPI (Message Passing Interface) [14] synchronous communication model to use Pthreads [1]. A thread of execution is a fork of a computer program into more concurrently running tasks. Those threads are executed independently but share memory and other resources. Typically, creation, destruction and inter-process communication using threads are faster than using processes. There exist more implementations of threads. In our case we have used Pthreads. Pthread is the *POSIX1003.1c* thread standard put out by the IEEE standards committee. This standard got the IEEE Standards Board approval in June 1995. In the MPI2 [14] specification there exists thread support inside of MPI. To achieve that only one thread can access shared memory we have used Pthread mutexes. We have used locking file mechanism to achieve exclusive access to shared file. The share file used in our code is accessed via network file system disallowing us to use *flock* or *fcntl* mechanism. For these reasons we have used only exclusive access and we have created lock file every time when thread attempts to read or write from the shared file. In Table2 and Table3 we report measured time spent in waiting for releasing mutexes and file locks including storing and loading of data.

We have used coarse-grained parallelization strategy where one colony is assigned to each thread. It means, that the number of parallel threads is equal to the number of colonies and each thread calculates solutions of all ants in a colony. All the time a better solution is found, it is stored in the shared memory within one node. The pseudoalgorithm can be formulated as follows:

```

1: Initialization;
2: For i=1; i<=It_out do:
    Reset pheromone matrix;

```



```

For j=1; j<=It_in do:
  For Ant=1; Ant<=n/2 do:
    Create Savings based Ant solution;
    Select elitist Ants;
    Update solution within one node if better solution is found;
    Update pheromone matrix if needed;
    Update Savings list if needed;
  End do j;
  Update solution between nodes in shared files if better;
  solution is found;
End do i;
3: Finalization;

```

Table 1: Measured problem instances, where  $n$  denotes the number of customers and  $Q$  denotes the vehicle capacity.

Instance	$n$	$Q$
C4	150	200
C5	199	200
G18	300	200
G19	360	200
G20	420	200

## 4 Computational results

In our experiments we have used test instances presented in Table 1. They are two larger instances generated by Christofides et al. [4] and three of the larger instances generated by Golden et al. [12]. For all presented experiments we have used the cluster<sup>1</sup>, University of Vienna, consisting of 72 SUN X4100 nodes with two 64-bit dual core processors, each. Therefore we could use 4 threads working over the common shared memory. For communication within one node we have used common shared memory and the communication between nodes was implemented via storing of exchanged data in shared files over the network. The threads are divided into groups, whereby the rank of them is dependent on the number of node cores. For communication between threads belonging to the same group, shared memory is used. For communication between nodes the file access is used. The communications are asynchronous and are realized after a better solution is found.

Reported results when not mentioned otherwise are average values gained over 30 runs for each instance. Denoting by  $n$  the number of customers, for all instances we have used these

---

<sup>1</sup>For details see: <http://luna.cs.univie.ac.at/aurora/description.htm>

Table 2: Calculated average results according to the number of threads of each measured instance, where  $Th$  denotes number of the threads used,  $V$  denotes calculated quality solution,  $t_c$  denotes the overall time spent by communication,  $t_r$  denotes the overall execution time (communication time included) and  $I_s$  denotes the iteration number where solution value has been stabilized

Instance	$Th$	$V$	$t_c[s]$	$t_r[s]$	$I_s$
C4	1	1070.73	1.96	52.10	78.97
	4	1064.48	9.37	53.87	91.75
	8	1056.99	10.71	52.53	174.75
	16	1052.99	13.37	52.42	232.50
	32	1048.08	21.06	54.37	240.61
C5	1	1377.18	1.68	143.24	95.28
	4	1365.42	6.63	141.90	110.62
	8	1352.98	7.38	132.16	254.9
	16	1346.58	8.24	127.93	282.14
	32	1337.69	10.94	126.73	314.18
G18	1	1090.35	1.53	504.58	112.97
	4	1081.58	5.67	505.18	146.87
	8	1074.12	10.81	471.68	290.35
	16	1066.00	6.17	451.48	341.46
	32	1061.18	7.69	436.50	368.32
G19	1	1501.36	1.57	1105.70	150.14
	4	1489.89	5.62	1105.50	156.82
	8	1479.22	5.65	1001.51	288.80
	16	1472.19	5.80	990.94	316.09
	32	1461.76	8.87	948.26	349.30
G20	1	2009.37	1.69	2058.67	161.45
	4	1994.98	5.78	2070.65	181.86
	8	1971.36	5.83	1849.31	303.98
	16	1956.45	5.90	1721.34	372.63
	32	1948.92	9.59	1632.35	355.42

Table 3: Calculated average results of C4 instance within one node without file access according to number of threads, where  $Th$  denotes number of used threads,  $V$  denotes the calculated quality solution,  $t_c$  denotes overall time spent by communication and  $I_s$  denotes iteration number where solution value has been stabilized

$Th$	$V$	$t_c[\text{ms}]$	$I_s$
1	1068.09	0.25	74.1
2	1063.27	0.66	76.32
3	1063.24	1.29	121.80
4	1063.76	1.30	190.30

configurations: We have used  $\lfloor n/2 \rfloor$  artificial ants for each thread,  $\alpha = \beta = 5$  and  $\sigma = 3$  elitist ants, the evaporation rate  $\rho = 0.95$ , and the neighborhood size  $\lfloor n/4 \rfloor$ . We let the algorithm run for  $It_{out} = 20$  outer iterations for each problem instance, and  $It_{in} = 20$  inner loops. After all ants in a colony found their solutions, the best  $\sigma$  solutions were chosen. They were compared with the best solutions found so far and saved in the common shared memory. If the last generated solution was not better, it was not saved and the shared memory stayed unchanged. The pheromone matrix was updated after a better solution was found. After each of the 20 outer iterations, the gained solutions between nodes were compared and updated. When comparing communication times in Table 2 and Table 3, we can see that using internode shared files dramatically increases the communication time. The Table 2 illustrates the fact that the time spent by communication increases with using more cluster nodes. The communication time increases only slightly with increasing of number of threads in the frame of one node. This seems to be caused by fact that communication is done only if better solution is found within node (4 threads). In the case of using just one thread for calculation, the file access is performed each time better solution is found. From Table 2 we can see that the solution quality is increasing with the number of colonies, because each colony corresponds to a separate thread. The reason is that more artificial ants and therefore more routes are generated. This fact supports the advantages of executing the parallel algorithm within the multicore processor environment. It is interesting that the increased number of threads and so the colonies, increases also the number of iteration where the best quality solution has been found.

## 5 Conclusions

We have presented parallel Posix threads based implementation of Ant Colony Optimization method for solving the Vehicle Routing Problem. The algorithm is based on a homogeneous multi-colony approach and used asynchronous communication in finding solutions. We have showed that the quality of solution is improved in dependence on the number of colonies, whereby each colony is assigned to a separate thread.

In our future work we would like to study various asynchronous algorithms for ACO and test

them on multicore architectures with more cores.

### Acknowledgement

The authors would like to thank Prof. Siegfried Benkner, Institute of Scientific Computing, University of Vienna, for the possibility to use the parallel computer cluster LUNA for running their experiments.

### References

- [1] ANDREWS, G.R.: *Foundations of Multithreaded, Parallel, and Distributed Programming*. Addison-Wesley, Reading, MA, 2000.  
<http://www.cs.arizona.edu/people/greg/mpdbook>
- [2] BELL, J.E., McMULLEN, P.R.: *Ant colony optimization techniques for vehicle routing problem*. Advanced Engineering Informatics, 18, pp.41–48, 2004.
- [3] CAPELLO, F., ETIEMBLE, D.: *MPI versus MPI+OpenMP on the IBM SP for the NAS benchmarks*. In Proceedings of SuperComputing 2000. IEEE Computer Society, 2000.
- [4] CHRISTOFIDES, N., MINGOZZI, A., TOTH, P.: *The Vehicle Routing Problem*. In CHRISTOFIDES, N., MINGOZZI, A., TOTH, P., SANDI, C., (Eds.): *Combinatorial Optimization*, Wiley, Chichester, 1979.
- [5] CIONI, L.: *Some Strategies for Parallelizing Ant Systems*. Doctoral Course Parallel computing in Combinatorial Optimization, 2005.
- [6] CRAINIC, T. G.: *Parallel Solution Methods for Vehicle Routing Problems*. CIRRELT-2007-28, 2007.
- [7] DOERNER, K. F., HARTL, R.F., BENKNER, S., LUCKA, M.: *Cooperative Savings based Ant Colony Optimization - Multiple Search and Decomposition Approaches*. Parallel Processing Letters, 16(3), pp.351–369, 2006.
- [8] DOERNER, K.F., HARTL, R.F., REINMANN, M.: *Ants solve time constrained pickup and delivery problems with full truckloads*. In B. Fleishmann, R. Lasch, U. Derigs, W. Domschke and U. Reider, eds: *Operations Research Proceedings 2000*, Springer, Berlin, 395–400, 2001.
- [9] DORIGO, M., GAMBARDELLA, L. M.: *Ant Colony System: A cooperative learning approach to the Travelling Salesman Problem*. IEEE Transactions on Evolutionary Computation, 1(1), pp.53–66, 1997.
- [10] GAMBARDELLA, L.M., TAILLARD, É., AGAZZI, G.: *MACS-VRPTW: a multiple ant colony system for vehicle routing problems with time windows*. In D. Corne, M. Dorigo and F. Glover, eds: *New ideas in optimization*, McGraw-Hill, London, pp.63–76, 1999.
- [11] GENDREAU, M., POTVIN, J. Y., BRÄYSY, O., HASLE, G., LOKKETANGEN, A.: *Metaheuristics for Vehicle Routing Problem and its Extensions : A Categorized Bibliography*. CIRRELT-2007-27, 2007.
- [12] GOLDEN, B. L., WASIL, E. A., KELLEY, J. P., CHAO, K. M.: *The impact of metaheuristics on solving the vehicle routing problem: algorithms, problem sets, and computational results*. In T. G. Crainic, G. Laporte, eds: *Fleet management and logistics*, Norwell: Kluwer, 1998.

- [13] MAZZEO, S., LOISEAU, I.: *An ant colony algorithm for capacitated vehicle routing*. Electronic Notes in Discrete Mathematics, 18, pp.181–186, 2004.
- [14] Message Passing Interface Forum. *MPI: A Message-Passing Interface Standard*. Vers. 1.1, June 1995. MPI-2: Extensions to the Message-Passing Interface, 1997.
- [15] <http://luna.cs.univie.ac.at>

#### Current address

**Maria Lucka, doc. RNDr. CSc.**

Faculty of Education, University of Trnava, 918 43 Trnava, Priemyselna 8,  
e-mail: [mlucka@truni.sk](mailto:mlucka@truni.sk)

**Stanislav Piecka, Ing.**

Faculty of Controlling and Informatics, University of Zilina, 010 26 Zilina, Univerzitna 8215/1,  
e-mail: [piecka@gmail.com](mailto:piecka@gmail.com)



# ON ONE ORTHOGONAL TRANSFORM APPLIED ON A SYSTEM OF ORTHOGONAL POLYNOMIALS IN TWO VARIABLES

MARČOKOVÁ Mariana, (SK), GULDAN Vladimír, (SK)

**Abstract.** We use an orthogonal transform by goniometric functions to the partial differential equations for Jacobi orthogonal polynomials in two variables taken as products of classical Jacobi polynomials in one variable. The results are specified for Legendre polynomials in two variables and functions associated with them which have wide applications.

**Key words and phrases.** orthogonal polynomial, orthogonal polynomial in two variables, second order partial differential equation, Jacobi polynomial, Legendre polynomial.

*Mathematics Subject Classification.* Primary 33C30, Secondary 35G05.

## 1 Introduction

V. Jarník in his book [4] introduces the term "orthogonal transform" that is often used in mathematical modelling in natural and technical sciences (cf. [1-2] and [6]). In the Jarník's book we can find some examples of orthogonal transforms, e.g. transform by polar coordinates, transform by spherical coordinates, etc.

In the present paper we use orthogonal transform by goniometric functions. We apply them to the partial differential equations that are satisfied by classical Jacobi orthogonal polynomials in two variables.

It is well known that the classical Jacobi polynomials  $\{J_n(x)\}_{n=0}^{\infty}$  are orthogonal in the interval  $< -1, 1 >$  with respect to the weight function

$$J(x) = (1-x)^{\alpha}(1+x)^{\beta}.$$

In this contribution we consider their two-dimensional pendants defined as the products  $J_n(x)J_m(y)$ ,  $n = 0, 1, 2, \dots$ ,  $m = 0, 1, 2, \dots$  which are orthogonal in the square region

$$R = \{(x, y); -1 < x < 1, -1 < y < 1\}$$

with respect to the weight function

$$J(x)J(y) = (1-x)^\alpha(1+x)^\beta(1-y)^\alpha(1+y)^\beta.$$

It means that

$$\iint_R J_n(x)J_m(y)J_k(x)J_l(y)(1-x)^\alpha(1+x)^\beta(1-y)^\alpha(1+y)^\beta dx dy = 0$$

for  $(n, m) \neq (k, l)$ ,  $n = 0, 1, 2, \dots$ ,  $m = 0, 1, 2, \dots$ ,  $k = 0, 1, 2, \dots$ ,  $l = 0, 1, 2, \dots$ .

## 2 Orthogonal transform of partial differential equations satisfied by $J_n(x)J_m(y)$

It is well known (cf. [5, 7-8]) that

$$(1) \quad D_1^{\alpha, \beta} J_n(x)J_m(y) = [-n(n + \alpha + \beta + 1) - m(m + \alpha + \beta + 1)]J_n(x)J_m(y)$$

where

$$D_1^{\alpha, \beta} = (1-x^2)\frac{\partial^2}{\partial x^2} + (1-y^2)\frac{\partial^2}{\partial y^2} + [\beta - \alpha - (\alpha + \beta + 2)x]\frac{\partial}{\partial x} + [\beta - \alpha - (\alpha + \beta + 2)y]\frac{\partial}{\partial y}$$

is the second order partial differential operator defined for  $\alpha > -1$ ,  $\beta > -1$ .

Let us transform the operator  $D_1^{\alpha, \beta}$  into another form by the orthogonal transform

$$(2) \quad x = \cos u, \quad y = \sin v,$$

i.e.,

$$u(x) = \arccos x, \quad v(y) = \arcsin y.$$

It means that the region  $R = \{(x, y); -1 < x < 1, -1 < y < 1\}$  is transformed to the region

$$U = \left\{ (u, v); 0 < u < \pi, -\frac{\pi}{2} < v < \frac{\pi}{2} \right\}.$$

Then we have

$$\frac{du}{dx} = -\frac{1}{\sqrt{1-x^2}}, \quad \frac{dv}{dy} = \frac{1}{\sqrt{1-y^2}}$$

and

$$\begin{aligned} \frac{\partial}{\partial x} &= \frac{\partial}{\partial u} \frac{du}{dx} = -\frac{\partial}{\partial u} \frac{1}{\sqrt{1-x^2}}, \quad \frac{\partial}{\partial y} = \frac{\partial}{\partial v} \frac{dv}{dy} = \frac{\partial}{\partial v} \frac{1}{\sqrt{1-y^2}}, \\ \frac{\partial^2}{\partial x^2} &= \frac{\partial^2}{\partial u^2} \frac{1}{1-x^2} - \frac{\partial}{\partial u} \frac{x}{(1-x^2)^{\frac{3}{2}}}, \quad \frac{\partial^2}{\partial y^2} = \frac{\partial^2}{\partial v^2} \frac{1}{1-y^2} + \frac{\partial}{\partial v} \frac{y}{(1-y^2)^{\frac{3}{2}}}. \end{aligned}$$



So the operator  $D_1^{\alpha,\beta}$  is transformed to the operator

$$(3) \quad D_2^{\alpha,\beta} = \frac{\partial^2}{\partial u^2} + \frac{\partial^2}{\partial v^2} + [(\alpha - \beta) \csc u + (\alpha + \beta + 1) \cot u] \frac{\partial}{\partial u} + [(\beta - \alpha) \sec v - (\alpha + \beta + 1) \tan v] \frac{\partial}{\partial v} .$$

It is obvious that the equations

$$D_2^{\alpha,\beta} f(u, v) = [-n(n + \alpha + \beta + 1) - m(m + \alpha + \beta + 1)] f(u, v)$$

are satisfied by the products  $J_n(\cos u) J_m(\sin v)$ ,  $n = 0, 1, 2, \dots, m = 0, 1, 2, \dots$ .

In (3) we denote

$$\omega_1(u) = (\alpha - \beta) \csc u + (\alpha + \beta + 1) \cot u$$

and

$$\omega_2(v) = (\beta - \alpha) \sec v - (\alpha + \beta + 1) \tan v .$$

Further, denote

$$q(u, v) = \sqrt{\sin u \cos v J(\cos u) J(\sin v)} .$$

We express its natural logarithm

$$\begin{aligned} \ln q(u, v) &= \frac{1}{2} [\ln \sin u + \ln \cos v] + \frac{\alpha}{2} [\ln(1 - \cos u) + \ln(1 - \sin v)] + \\ &+ \frac{\beta}{2} [\ln(1 + \cos u) + \ln(1 + \sin v)] . \end{aligned}$$

Then

$$(4) \quad \frac{\partial \ln q(u, v)}{\partial u} = \frac{\frac{\partial q(u, v)}{\partial u}}{q(u, v)} = \frac{1}{2} \omega_1(u)$$

and

$$(5) \quad \frac{\partial \ln q(u, v)}{\partial v} = \frac{\frac{\partial q(u, v)}{\partial v}}{q(u, v)} = \frac{1}{2} \omega_2(v) .$$

So

$$\frac{\frac{\partial q(u, v)}{\partial u} + \frac{\partial q(u, v)}{\partial v}}{q(u, v)} = \frac{1}{2} [\omega_1(u) + \omega_2(v)]$$

and integrating (4) and (5) from 0 to  $u$  and  $v$ , respectively, where  $u \in (0, \pi)$ ,  $v \in (-\frac{\pi}{2}, \frac{\pi}{2})$  we have

$$\ln q(u, v) = \frac{1}{2} \left( \int_0^u \omega_1(t) dt + \int_0^v \omega_2(t) dt \right)$$

and

$$q(u, v) = \exp \left[ \frac{1}{2} \left( \int_0^u \omega_1(t) dt + \int_0^v \omega_2(t) dt \right) \right] .$$

For  $n = 0, 1, 2, \dots$ ,  $m = 0, 1, 2, \dots$  we take the functions

$$(6) \quad \begin{aligned} q_{n,m}(u, v) &= J_n(\cos u) J_m(\sin v) q(u, v) = \\ &= J_n(\cos u) \exp \left[ \frac{1}{2} \int_0^u \omega_1(t) dt \right] J_m(\sin v) \exp \left[ \frac{1}{2} \int_0^v \omega_2(t) dt \right]. \end{aligned}$$

The functions (6) satisfy the differential equations

$$(7) \quad \begin{aligned} &\frac{\partial^2 \varphi_{n,m}(u, v)}{\partial u^2} + \frac{\partial^2 \varphi_{n,m}(u, v)}{\partial v^2} = \\ &= \left[ \frac{\omega_1^2(u) + \omega_2^2(v)}{4} + \frac{\omega_1'(u) + \omega_2'(v)}{2} - n(n + \alpha + \beta + 1) - m(m + \alpha + \beta + 1) \right] \times \\ &\quad \times \varphi_{n,m}(u, v) \end{aligned}$$

in the region  $U$  (cf.[3, p.123]). So, the differential equations (7) are of the type

$$\Delta \varphi_{n,m}(u, v) = \lambda_{n,m}(u, v) \varphi_{n,m}(u, v)$$

where  $\Delta$  is the Laplace operator of  $u$  and  $v$  and

$$\lambda_{n,m}(u, v) = \frac{\omega_1^2(u) + \omega_2^2(v)}{4} + \frac{\omega_1'(u) + \omega_2'(v)}{2} - n(n + \alpha + \beta + 1) - m(m + \alpha + \beta + 1).$$

### 3 Orthogonality of the functions $q_{n,m}(u, v)$

Suppose that  $(n, m) \neq (k, l)$ . Then

$$\begin{aligned} &\iint_U q_{n,m}(u, v) q_{k,l}(u, v) du dv = \\ &= \iint_U J_n(\cos u) J_m(\sin v) J_k(\cos u) J_l(\sin v) q^2(u, v) du dv = \\ &= \iint_U J_n(\cos u) J_m(\sin v) J_k(\cos u) J_l(\sin v) \sin u \cos v J(\cos u) J(\sin v) du dv = \\ &= \int_0^\pi J_n(\cos u) J_k(\cos u) \sin u J(\cos u) du \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} J_m(\sin v) J_l(\sin v) \cos v J(\sin v) dv = \\ &= \int_{-1}^1 J_n(x) J_k(x) J(x) dx \int_{-1}^1 J_m(y) J_l(y) J(y) dy = 0 \end{aligned}$$

because at least one of the last integrals is equal to zero.

Similarly, if the systems  $\{J_n(x)\}_{n=0}^{\infty}$  and  $\{J_m(y)\}_{m=0}^{\infty}$  are orthonormal in the interval  $< -1, 1 >$  we have

$$\begin{aligned} \iint_U q_{n,m}^2(u, v) du dv &= \iint_U J_n^2(\cos u) J_m^2(\sin v) q^2(u, v) du dv = \\ &= \iint_U J_n^2(\cos u) J_m^2(\sin v) \sin u \cos v J(\cos u) J(\sin v) du dv = \\ &= \int_0^{\pi} J_n^2(\cos u) \sin u J(\cos u) du \int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} J_m^2(\sin v) \cos v J(\sin v) dv = \\ &= \int_{-1}^1 J_n^2(x) J(x) dx \int_{-1}^1 J_m^2(y) J(y) dy = 1. \end{aligned}$$

Thus we proved that the functions  $q_{n,m}(u, v)$  given by (6) are orthogonal (orthonormal) functions in two variables on  $U$  with respect to the weight function  $w(u, v) = 1$ .

#### 4 Legendre polynomials and Legendre associated functions in two variables

If in (1) we have  $\alpha = \beta = 0$ , then we get

$$D_1^{0,0} = (1 - x^2) \frac{\partial^2}{\partial x^2} + (1 - y^2) \frac{\partial^2}{\partial y^2} - 2x \frac{\partial}{\partial x} - 2y \frac{\partial}{\partial y}$$

and

$$D_1^{0,0} L_n(x) L_m(y) = [-n(n+1) - m(m+1)] L_n(x) L_m(y)$$

where  $\{L_n(x)\}_{n=0}^{\infty}$  and  $\{L_m(y)\}_{m=0}^{\infty}$  are the Legendre polynomials which are orthogonal for  $-1 < x < 1$  and  $-1 < y < 1$  with respect to the weight functions  $L(x) = 1$  and  $L(y) = 1$  (cf. [9]).

From there

$$D_2^{0,0} = \frac{\partial^2}{\partial u^2} + \frac{\partial^2}{\partial v^2} + \cot u \frac{\partial}{\partial u} - \tan v \frac{\partial}{\partial v}$$

and

$$D_2^{0,0} L_n(\cos u) L_m(\sin v) = [-n(n+1) - m(m+1)] L_n(\cos u) L_m(\sin v).$$

Then

$$q(u, v) = \sqrt{\sin u \cos v}$$

and the equations

$$\Delta \varphi_{n,m}(u, v) =$$

$$= \left[ \frac{\cot^2 u + \tan^2 v}{4} + \frac{1}{2} \left( \frac{1}{\sin^2 u} - \frac{1}{\cos^2 v} \right) - n(n+1) - m(m+1) \right] \varphi_{n,m}(u, v)$$

are satisfied by the functions

$$q_{n,m}(u, v) = L_n(\cos u) L_m(\sin v) \sqrt{\sin u \cos v}$$

which are orthogonal on the square region  $U$  with respect to the weight function  $w(u, v) = 1$ .

It can be concluded that the functions  $q_{n,m}(u, v)$  taken as the solutions of the partial differential equations (1) in the case of products of Legendre polynomials  $L_n(x)L_m(y)$  are after orthogonal transform (2) in the role of associated Legendre functions in two variables which are orthogonal, too. They have wide applications in physics and other natural and technical sciences.

### Acknowledgement

The paper was supported by the Grant Agency VEGA of Ministry of Education of Slovak Republic through the grant no. 1/0867/08 with the title "Properties of Orthogonal Systems Applied in Natural and Engineering Sciences".

### References

- [1] DOBRUCKÝ, B., MARČOKOVÁ, M., POKORNÝ, M., ŠUL, R.: *Using Orthogonal and Discrete Transform for Single - Phase PES Transients - a New Approach*. Proceedings of the 27th IASTED International Conference Modelling, Identification, and Control, February 11 - 13, 2008, Innsbruck, Austria, 60 - 65, 2008.
- [2] DOBRUCKÝ, B., ŠUL, R., BEŇOVÁ, M.: *Mathematical Modelling of Two-Stage Converter using Complex Conjugated Magnitudes and Orthogonal Park-Clarke Transformation Methods*. Submitted to Aplimat 2009.
- [3] GREGUŠ, M., ŠVEC, M., ŠEDA, V.: *Ordinary Differential Equations*. Alfa, Bratislava, 1985 (in Slovak).
- [4] JARNÍK, V.: *Calculus (II)*. Academia, Praha, 1984 (in Czech).
- [5] KOORNWINDER, T. H.: *Orthogonal Polynomials in Two Variables which are Eigenfunctions of Two Algebraically Independent Partial Differential Operators*. Indag. Math., Vol. 36, pp. 48-58, 1974.
- [6] MAMRILLA, D.: *On the Systems of First Order Quasi-linear Differential Equations*. Es-sox, Prešov, 2004 (in Slovak).
- [7] MARČOKOVÁ, M.: *Approximation of Functions in Two Variables by Cesáro Means of Fourier-Jacobi Sums*. Proceedings of International Scientific Conference of Mathematics, Žilina, 1998, 161-165, 1999.
- [8] MARČOKOVÁ, M.: *Second Order Partial Differential Equations for Some Orthogonal Polynomials in Two Variables*. Studies of University of Žilina, math. - phys. series 13, 127-132, 2001.
- [9] SUJETIN, P. K.: *Classical Orthogonal Polynomials*. Nauka, Moskva, 1979 (in Russian).

**Current address**

**Mariana Marčoková, doc., RNDr., CSc.**

Faculty of Science, University of Žilina, Univerzitná 1, 010 26 Žilina,

e-mail: mariana.marcokova@fpv.uniza.sk

**Vladimír Guldán, RNDr.**

Gymnázium Hlinská, 010 01 Žilina

e-mail: vladimir.guldan@post.sk



## AN INTROSPECT OF SIMULATION OF NONDETERMINISTIC TURING MACHINE WITH A REAL-ANALYTIC FUNCTION

PIEKARZ Monika, (PL)

**Abstract.** In this paper some modification of a real-analytic simulation of nondeterministic Turing machine is given by means of finite state automata and Collatz function. This is the modification of the simulation presented in [7] and based on the paper of Koiran and Moore [4]. The simulation presented here is less complex than in [7].

**Key words and phrases.** nondeterministic Turing machine, simulation.

*Mathematics Subject Classification.* Primary 68Q10

### 1 Introduction and motivation

The Turing machine is the most useful model to point out complexity classes of some problems [9, 11]. Recently several types of Turing machines, which can solve some problems undecidable by the classical Turing machines, have been introduced [3, 5, 10]. Notwithstanding in complexity theory the main role plays non-realistic model called nondeterministic Turing machine (*NTM*) which expands a deterministic version.

In this paper we'll present some simulation of nondeterministic Turing machine by real-analytic function built from elementary functions. Systems of functions in which such a simulation is possible are able to help us to find solutions of classical (un)decidability problems which could be analyzed in models with possible simulations. Also a simplicity of simulations can be considered from the mathematical point of view. Various authors have independently shown that finite-dimensional piecewise-linear maps and flows can simulate Turing machines. Such methods have been presented in various papers [4, 12, 13] sometimes with additional requirements for simulating functions.

In [7] a simulating real-analytic function  $f$  which for some number coding input configuration of  $NTM$  give us after  $t'$  iterations some number coding output configuration of  $NTM$  was presented,  $t'$  has an exponential growth respect to time-step of the  $NTM$ . Simulation presented here is the modification of simulation presented in [7]. Simulating function which we will present now is similar form as in [7] but it is constructed from  $210k$  terms instead  $2310k$  terms where  $k$  is number of states of  $NTM$ .

## 2 Basic definition

We make now some introduction to present the main result.

Formally, a nondeterministic Turing machine ( $NTM$ ) of one tape is defined as a 5-tuple  $(Q, \Sigma, \delta, q_0, q_f)$ , where  $Q$  is a nonempty finite set (of  $k = |Q|$  states),  $\Sigma$  is not empty finite set of a tape alphabet (of  $m = |\Sigma|$  symbols),  $q_0 \in Q$  is the initial state,  $q_f \in Q$  is the accepting state, and  $\delta : Q \times \Sigma \rightarrow \wp(Q \times \Sigma \times \{\leftarrow, \rightarrow\})$  is a total function called a transition function of  $M$ , where, for any set  $A$ ,  $\wp(A)$  denotes the power set of  $A$ .

An input  $w$  is accepted by a nondeterministic Turing machine if and only if there exists a computation of this  $NTM$  on  $w$  ending in the accepting state. In this paper we will state that all symbols from  $Q$  and  $\Sigma$  are coded as a sequence of natural numbers, so  $Q = \{0, 1, \dots, k-1\}$  and  $\Sigma = \{0, 1, \dots, m-1\}$  where 0 is the code of the empty symbol which fills the whole tape besides the finite number of tapes cells which are used by  $NTM$  during the computation. Moreover, we don't allow  $NTM$  with more then two next possible movements. More possibilities are useless because the computation is only a constant factor longer when we bound to two possible choices in one step of the computation. Thus, without loss of generality, if we require in our model that for every pair  $(q, a)$  from  $Q \times \Sigma$ ,  $|\delta(q, a)| \leq 2$ , specifying  $\hat{\delta} : Q \times \Sigma \times \{0, 1\} \rightarrow Q \times \Sigma \times \{\leftarrow, \rightarrow\}$ , then bits 0 and 1 prescribe at most two branches of the computation in each state for each symbol being read.

## 3 Simulations of nondeterministic Turing machine by $FSA$

We present a modification of the simulation given by Koiran and Moore in [4]. We have presented similar simulation in [7]. However, we consider here the simulation with less complex finite automaton used as a help to find a simulation function. So the simulation function which we receive is less complex, too.

Now we shortly recall an idea of the simulation. In the construction we will use a finite state automaton  $FSA$  which has a finite number of counters and can increment or decrement the counters or check their equality to 0 ([6]).

Basic difference between this simulation and the simulation proposed in [7] is that the current  $FSA$  uses four counters instead five, they are:  $L$  - the code of the left part of a tape,  $R$  - the code of the right part of a tape ( $R$  includes the tape symbol  $a_0$  at the head current location),  $G$  - the code of guess,  $W$  is a working counter. Specifically, if our  $NTM$  has  $m$  tape symbols,  $L = \sum_{i=1}^{\infty} m^{i-1} a_{-i}$ ,  $R = \sum_{i=0}^{\infty} m^i a_i$  and  $G = \sum_{i=0}^{t-1} 2^i g_i$  where  $t$  is the number of steps of  $NTM$  computation and  $g_i \in \{0, 1\}$ . In [7] we had the fifth counter  $S$  which contained the current state of  $NTM$ . Now we'll remember the current state of  $NTM$  in states of  $FSA$ .



Every particular variant of one step of the transition function  $\hat{\delta}$  is described by three blocks of nodes of  $FSA$ . So this time a general description of the  $FSA$  simulating the behavior of the  $NTM$  is given on Fig. 1.

Let us explain the meaning of particular elements of  $FSA$ . Block I (Fig. 2) is constructed for the purpose of checking a current symbol under the head of  $NTM$ .

The counters on the output of Blok I are equal to  $W = \lfloor R/m \rfloor$ ,  $R = 0$ ; moreover we know the symbol  $a$  which is actually under the head, where  $\lfloor \rfloor$  denotes cut to the integer part.

In the case when state  $s$  is a final state of  $NTM$  Blok I is formed as on Fig. 3 and  $FSA$  finishes its computation but in a rest case it gets to the Blok III (Fig. 4), where our model has to determine the next bit of a guess  $p$ . Output values of counters of this blok are equal  $G = \lfloor G/2 \rfloor$ ,  $W = 0$ .

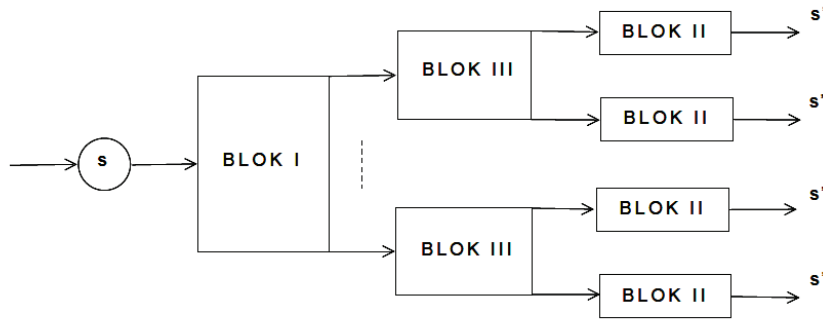


Figure 1:  $FSA$  simulating the behavior of  $NTM$ .

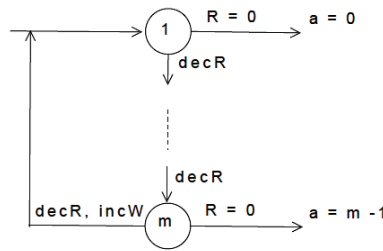


Figure 2: Blok I. Checking the head of  $NTM$ .

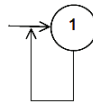


Figure 3: Blok I in the case when  $s$  is the final state of  $NTM$ .

Now we must fork the remaining part of the scheme. In the case when transition function is  $\hat{\delta}(s, a, p) = (s', a', \rightarrow)$ , i.e. the move in the right direction, we finish with Block II (Fig. 5)

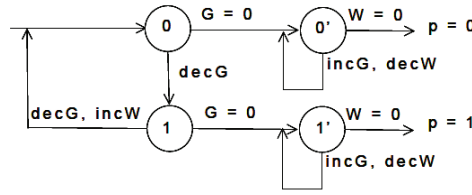


Figure 4: Blok III. Checking one bit of the guess of  $NTM$ .

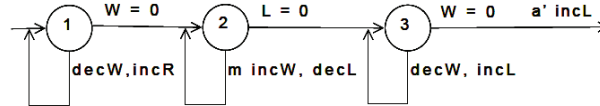


Figure 5: Blok II. The right direction move, changing the symbol under the head.

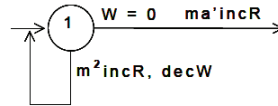


Figure 6: Blok IIa. Changing the symbol under the head.

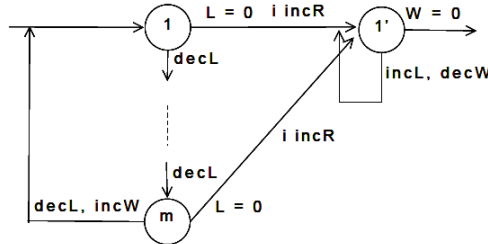


Figure 7: Blok IIb. The left direction move, checking the symbol on the left to the head and changing the counters  $FSA$  respectively to the left move of the head.

which ends with counters equal respectively ( $a'$  is the symbol given by the transition function for symbol  $a$  and state  $s$ ):  $R = \lfloor R/m \rfloor$ ,  $L = mL + a'$ ,  $W = 0$ .

However, the continuation of Block III in the case  $\hat{\delta}(s, a, p) = (s', a', \leftarrow)$  is more complicated and it consists of two parts: Block IIa (Fig. 6) and Block IIb (Fig. 7).

At the end of Block IIa we have counters equal respectively: ( $a'$  is the symbol given by the transition function for symbol  $a$  and state  $s$ )  $R = m^2R + ma'$ ,  $W = 0$  and we are ready to the left shift, so we can finish with Block IIb which checks the symbol on the left to the head and add it to the  $R$ . In this case we get the following results:  $R = m^2\lfloor R/m \rfloor + ma' + i$ ,  $L = \lfloor L/m \rfloor$ ,  $W = 0$  where  $i$  is the checked symbol on the left to the head.

After Block II in both cases our  $FSA$  changes its state into the state corresponding with a new state  $s'$  of  $NTM$ .

The simulation of each step of  $NTM$  needs  $O(L, R, G)$  steps of  $FSA$ ,  $L$  and  $R$  are size  $O(m^l)$  where  $l$  is the tape length used by  $NTM$ ,  $G$  is size  $2^t$  and  $t$  is the time which takes

the computation of nondeterministic Turing machine. So we can simulate it by *FSA* with 4 counters in time  $O(t(m^l + 2^t))$ .

#### 4 Simulation the FSA by real-analytic function of one variable

We now show how to simulate such *FSA* with a real-analytic function, in this case the Collatz function. Here, our method is the same as in [4]. Let  $f(x) = a_i x + b_i$ , for  $x \equiv i \pmod{u}$  for some base  $u$  and constants  $a_i, b_i$ , where  $0 \leq i < u$ . We will call any of such  $f$  a Collatz function. We define  $x = 2^L 3^R 5^G 7^W k + s$  where  $k$  is the number of states and  $s$  is the current state of *FSA*. Clearly all of our operations can be carried out on  $x$ . For instance, to decrement  $W$ , increment  $R$   $m^2$  times, we write

$$(decW, m^2 incR) : f(x) = (3^{m^2}/7)(x - s) + s' \text{ so } a = 3^{m^2}/7, \text{ and } b = s' - (3^{m^2}/7)s.$$

Owing to the fact that our simulating *FSA* uses only four counters, we can test for zero on all our counters in terms of  $x \pmod{210k}$  instead in terms  $x \pmod{2310k}$  like in [7]. We use for them the special function  $h_{210k}$  with the following property  $h_{210k}(x - i)$  is equal to 1 iff  $x \equiv i \pmod{210k}$ . The function

$$h_u(x) = \left( \frac{\sin \pi x}{u \sin \frac{\pi x}{u}} \right) = \begin{cases} 1 & x \pmod{u} = 0 \\ 0 & x \pmod{u} \neq 0 \end{cases},$$

with  $u = 210k$  is suitable for the purpose. All important cases are given in Table 1.

Only  $104k$  terms are actually needed, since the other  $106k$  possibilities  $i = x \pmod{210}$  never happen. Then we have the one-dimensional simulation of nondeterministic Turing machines given by the real-analytic function  $f(x) = \sum_{i=0}^{p-1} h_{210k}(x - i)(a_i x + b_i)$ <sup>1</sup>.

If we now put  $s_0$  as a start state of *NTM*,  $g$  as a natural number coding guess of *NTM* in particular steps and input of *NTM* on the tape on the left of the head than  $x = 2^w 5^g + s_0$  is an equivalent value to the input  $w$  of *NTM*. So we can form the following theorem:

**Theorem 4.1** *For any nondeterministic Turing machine  $M$  with  $m$  tape symbols and any input  $w$ , and any guess  $g$ , there is a real-analytic function  $f$  of one variable and constants  $k$  (number of states of *NTM*) and  $s_0$ , such that  $M$  halts after  $t$  time-steps with a result  $y$  iff there exists  $t' \in \mathbb{N}$  such that  $f^{t'}(2^w 5^g + s_0) = 2^{y_1} 3^{y_2} 5^{g_1} k + s_1$  where  $y = y_1 \times m^{|y_2|} + y_2$  and  $s_1$  is a final state of  $M$  and  $t' = O(m^t + 2^t)$ .*

#### References

- [1] E. ASARIN, O. MALER and A. PNUELI, *Reachability analysis of dynamical systems with piecewise-constant derivatives*, Theor. Comp. Sci., 138: 35-66, 1995.
- [2] J.H. CONWAY, *Unpredictable iteration*, Proc. 1972 Number Theory University of Colorado, 49- 52, 1972.

<sup>1</sup>By real-analytic function we mean function which is sum of convergent series.

- [3] B. J. COPELAND, *Even Turing machines can compute uncomputable functions*, In: C. S. Calude, J. Casti, and M. J. Dinnen (eds), *Unconventional Models of Computation*, Springer-Verlag, 1998.
- [4] P. KOIRAN and C. MOORE, *Closed-form analytic maps in one and two dimensions can simulate universal Turing machines*, *Theoretical Computer Science*, 210(1): 217-223, 1999.
- [5] C. LOKHORST Gert-Jan, *Hypercomputation*, *Department of Rhilosophy*, University of Helsinki, 2001.
- [6] M. MINSKY, *Computation: Finite and Infinite Machines*, Prentice-Hall, 1967.
- [7] J. MYCKA, M. PIEKARZ, *Simulation of nondeterministic Turing machine with finite state automata*, *Proceedings of Aplimat 2005*, 323-328, 2005.
- [8] P. ODIFREDDI, *Classical recursion theory*, Elsevier, 1989.
- [9] P. ODIFREDDI, *Classical recursion theory II*. Elsevier, 1999.
- [10] T. ORD, *Hypercomputation: computing more than the Turing machine*, Honours Thesis, University of Melbourne, 2002.
- [11] C. H. PAPADIMITRIOU, *Computational Complexity*, Addison Wesley Longman a Person Education Company, 1994.
- [12] M. PIEKARZ, *Three simulations of Turing machines with use of real recursive function*, *Annales UMCS Informatica AI*, 2004.
- [13] H. SIEGELMANN and E. D. SONTAG, *On the computational power of neural nets*, *Journal of Comp. and Systems Sc.*, 50:132-150, 1995.

**University of Maria Curie-Skłodowska**

**Piekarz Monika, master**

pl. Marii Curie-Skłodowskiej 1, 20-031,

e-mail: mpiekarz@umcs.lublin.pl

Table 1: List of values  $x \bmod 210k$  relative to the values of counters of  $FSA$ 

rejestry	$i(x \bmod 210k)$
$L > 0, R > 0, G > 0, W > 0$	$s$
$L = 0, R > 0, G > 0, W > 0$	$s + 150k$
$L > 0, R = 0, G > 0, W > 0$	$s + 70k, s + 140k$
$L > 0, R > 0, G = 0, W > 0$	$s + 42k, s + 84k, s + 126k, s + 168k$
$L > 0, R > 0, G > 0, W = 0$	$s + 30k, s + 60k, s + 90k, s + 120k, s + 150k, s + 180k$
$L = 0, R = 0, G > 0, W > 0$	$s + 35k, s + 175k$
$L = 0, R > 0, G = 0, W > 0$	$s + 21k, s + 63k, s + 147k, s + 189k$
$L = 0, R > 0, G > 0, W = 0$	$s + 15k, s + 45k, s + 75k, s + 135k, s + 165k, s + 195k$
$L > 0, R = 0, G = 0, W > 0$	$s + 14k, s + 28k, s + 56k, s + 98k, s + 112k, s + 154k, s + 182k, s + 196k$
$L > 0, R = 0, G > 0, W = 0$	$s + 10k, s + 20k, s + 40k, s + 50k, s + 80k, s + 100k, s + 110k, s + 130k, s + 160k, s + 170k, s + 190k, s + 200k$
$L > 0, R > 0, G = 0, W = 0$	$s + 6k, s + 12k, s + 18k, s + 24k, s + 36k, s + 48k, s + 54k, s + 66k, s + 72k, s + 78k, s + 96k, s + 102k, s + 108k, s + 114k, s + 132k, s + 138k, s + 144k, s + 156k, s + 162k, s + 174k, s + 186k, s + 192k, s + 198k, s + 204k$
$L = 0, R = 0, G = 0, W > 0$	$s + 7k, s + 49k, s + 91k, s + 133k$
$L = 0, R = 0, G > 0, W = 0$	$s + 5k, s + 25k, s + 85k, s + 125k, s + 185k, s + 205k$
$L = 0, R > 0, G = 0, W = 0$	$s + 3k, s + 9k, s + 27k, s + 33k, s + 39k, s + 51k, s + 81k, s + 87k, s + 99k, s + 117k, s + 141k, s + 153k$
$L > 0, R = 0, G = 0, W = 0$	$s + 2k, s + 4k, s + 8k, s + 16k, s + 32k, s + 46k, s + 64k, s + 92k, s + 106k, s + 128k, s + 158k, s + 184k$

