## Interim Report          IR-01-066

## Risk Management: Modeling and Computer Applications

*Vyacheslav Maksimov (maksimov@imm.uran.ru)*
*Yuri Ermoliev (ermoliev@iiasa.ac.at)*
*Joanne Linnerooth-Bayer (bayer@iiasa.ac.at)*

Editors

**Approved by**

Arne B. Jernelöv (jernelov@iiasa.ac.at)
Acting Director

December 2001

# Risk Management:
# Modeling and
# Computer Applications

**Proceedings of IIASA Workshop**
**May 14-15, 2001**

**Vyacheslav Maksimov, Yuri Ermoliev and Joanne Linnerooth-Bayer**
**Editors**

IR-01-066 / December 2001

International Institute for Applied Systems Analysis

A-2361 Laxenburg, Austria

Russian National Member Organization, Moscow, Russia

# Foreword

This volume presents the first results obtained within the framework of a new research activity in the IIASA's Risk Modeling and Society (RMS) Project, which is supported by the Russian National Member Organization. The new activity is integrated with the RMS activities on insurance for natural hazards, which are presently focusing on methods to generate scenarios of catastrophic events, linking them to models of losses, and using numerical optimization techniques to improve the structure of insurance. Accordingly, a group of papers presented in this volume is devoted to different aspects of catastrophe modeling and insurance. Other important issues related to RMS's research interests, such as socio-economic and environmental aspects of risk management and advanced modeling techniques, are also discussed.

The papers collected in this volume were presented and discussed at the RMS-organized Workshop on Risk Management: Modeling and Computer Applications (IIASA, 14–15 May, 2001) in its four sessions: Seismic Models and Insurance; Risk Management and Modeling Techniques; Risk Management: Socio-Economic Aspects; and Risk Management: Optimization and Enviromental Aspects. The workshop agenda overlapped with research interests of IIASA's Environmentally Compatible Energy Strategies (ECS) Project and Dynamic Systems (DYN) Project, whose staff took part in the workshop. The DYN group helped in organizing the workshop.

Anielo Amendola, Marina Blizorukova, Yuri Ermoliev, Tatyana Ermolieva, Joanne Linnerooth-Bayer, Vyacheslav Maksimov, Valerii Rozenberg and Alexander Soloviev presented their papers in the session **Seismic Models and Insurance.**

The paper "Block models of lithosphere dynamics: approaches and solutions" by Alexander Soloviev and Vyacheslav Maksimov addresses the issue of catastrophe modeling. It is noted that the necessity of catastrophe modeling is punctuated by the increase of losses due to recent natural and anthropogenic hazards and by the lack of reliable observation data. For the case of earthquakes, models of the lithosphere dynamics constitute the basis of catastrophe modeling. The paper gives a brief overview of models of the lithosphere dynamics, with a focus on block models.

The paper "Risk insurance: generation of scenarios" by Marina Blizorukova, Yuri Ermoliev and Valerii Rozenberg (the authors acknowledge the assistance of Valerii Samosyuk) discusses the problems of providing insurance against natural hazards and helpful modeling tools. The focus is on methods for generating scenarios of earthquakes. It is conjectured that the generated scenarios may act as inputs to optimization algorithms aimed at improving the structure of insurance. The paper reviews relevant modeling approaches and their applicability given different classes of data.

The paper "Seismic risk management in the Toskany region: A stochastic optimization model" by Tatyana Ermolieva, Yuri Ermoliev and Anielo Amendola describes how a spatially dynamic stochastic optimization model that takes into account the complexities and interdependencies of

catastrophic risks can be customized to explicitly incorporate the geological characteristics of a region,the parameters of seismic hazards and the vulnerability of the built environment. The model is able to analyse multiple policy options for developing insurance in an equitable and fair manner, and the authors demonstrate what this means for insurance premiums and reserve funds. To analyze the stability of the system, the authors make use of non-differentiable stochastic optimization techniques combined with such measures of risk as Value-at-Risk (VaR) and the probability of bankruptcy.

Vladimir Kagramanian, Jozef Korbicz, Viktor Mourogov, Marcin Paprzycki, Valerii Rozenberg and Boris Digas were authors in the session **Risk Management and Modeling Techniques.**

The paper "The need of innovative fuel reactor and fuel cycle systems" by Vladimir Kagramanian and Viktor Mourogov has three sections: The first describes changes taking place in the market of nuclear power plants and the resulting needs for innovative reactors and fuel cycles. A second section outlines the range of innovative approaches in nuclear technologies, which have been identified so far, and the requirements for their successful implementation. Finally, the paper discusses the need for international cooperation in R&D, international initiatives that are already underway, and the role of the International Agency for Atomic Energy (IAEA).

The paper "Soft computing approaches in fault diagnosis system and risk management" by Jozef Korbicz discusses how neuron networks can be applied in different fields of science and technology. Such features of the neuron networks as simplicity in implementation, good approximation of built-up systems and the possibility of the convenient formation of equipment applications are emphasized. The authors note that increasing complexity of the examined objects and increasing requirements for efficiency and reliability of the applied analytical tools has stimulated the search for new solutions, in this context, research connected with integration of artificial neuron networks and other methods of artificial intelligence.

The paper "Medium structure modeling on parallel computers" by Marcin Paprzycki, Valerii Rozenberg, and Boris Digas deals with the application of methods for computer diagnostics to problems related to the assessment of natural and anthropogenic risk. The use of computational methods in the fields of plasma physics and geophysics are analyzed, a solution algorithm for a problem of reconstruction of the structure of a medium is described, and results of numerical experiments are discussed.

Marina Blizorukova, Andrei Maksimov, Oleg Nikonov and Andrei Shorikov presented their papers in the session on **Risk Management: Socio-Economic Aspects.**

The paper "Innovation: risk and economic safety of the Ural region" by Marina Blizorukova, Andrei Maksimov and Andrei Shorikov discusses problems concerned with the pace of the innovation processes in the Urals, a traditional industrial region of Russia. The authors highlight important aspects of risk management in the region and chraracterize the socio-economic aspects, as well as possible methods of risk assessment and innovation activity in small business.

The paper "Financial risk management: set-valued uncertainty modeling" by Oleg Nikonov deals with the problem of dynamic investment portfolio selection, which they treat with guaranteed control theory. A formalized setting and solution that combine the methods of this theory with the traditional mean-variance approach are discussed.

Sergei Aseev, Ger Klaassen, R. Alexander Roehrl and Alexander Tarasyev were the authors in the session **Risk Management: Optimization and Environmental Aspects.**

The paper "The Great Caspian Gas Pipeline Game" by Ger Klaassen, Alexander Roehrl and Alexander Tarasyev focuses on the problem of routing gas pipelines competing for the Turkey gas market. The authors propose a model of game-dynamic interactions between the pipeline projects. The model comprises four microeconomic levels of optimization: assessment of the market of potential innovations, selection of innovation scenarios, regulation of the future supply and optimization of the current investments. The projects interact through the macroeconomic price formation mechanism. The model is intended to serve as a macroeconomic tool for the analysis of the impacts of different investment policies for the construction of the pipelines. Of special interest for future research is the design of policies explicitly dealing with these types of risks.

The paper "Optimal control of dynamic system in presence of risky factors" by Sergei Aseev deals with the optimal control of dynamical systems whose state spaces contain domains of risk. Serious difficulties in the analysis of such problems arise due to the discontinuities with respect to the state variable, which may occur in the system's dynamics or in the cost functional. Two problems of optimal control with domains of risk are considered: a problem with state constraints and a problem of time-optimal crossing a given domain.

## Acknowledgments

# Contents

## PART I: SEISMIC MODELS AND INSURANCE

## PART II: RISK MANAGEMENT AND MODELING TECHNIQUES

# Participants

**Sergei Aseev**
Dynamic Systems
International Institute
For Applied Systems Analysis
A-2361 Laxenburg
AUSTRIA
E-mail: aseev@iiasa.ac.at

**Joanne Linnerooth-Bayer**
Risk, Modeling & Society
International Institute
for Applied Systems Analysis
A-2361 Laxenburg
AUSTRIA
E-mail: bayer@iiasa.ac.at

**Yuri Ermoliev**
IIASA Scholar
International Institute
for Applied Systems Analysis
A-2361 Laxenburg
AUSTRIA
E-mail: ermoliev@iiasa.ac.at

**Tatiana Ermolieva**
Social Security Reform
International Institute
for Applied Systems Analysis
A-2361 Laxenburg
AUSTRIA
E-mail: ermol@iiasa.ac.at

**Jozef Korbicz**
Technical University
of Zielona Gora
Ul. Podgorna 50
65-246 Zielona Gora
POLAND
E-mail: J.korbicz@irio.pz.zgora.pl

**Marina Blizorukova**
Institute of Mathematics and Mechanics
Urals Branch
Russian Academy of Sciences
Kovalevskaya str., 16
620066 Ekaterinburg
RUSSIA
E-mail: msb@imm.uran.ru

**Vyacheslav Maksimov**

Institute of Mathematics and Mechanics

Urals Branch

Russian Academy of Sciences

Kovalevskaya str., 16

620066 Ekaterinburg

RUSSIA

E-mail: maksimov@imm.uran.ru

**Oleg Nikonov**

Urals State Technical University

Mira str., 19

620002 Ekaterinburg

RUSSIA

E-mail: aspr@mail.ustu.ru

**Ger Klaassen**

Environmentally Compatible Energy

Strategies

International Institute

for Applied Systems Analysis

A-2361 Laxenburg

AUSTRIA

E-mail: klaassen@iiasa.ac.at

**Vladimir Kagramanian**

International Atomic Energy Agency

Wagramerstrasse 5

A-1400 Vienna

AUSTRIA

E-mail: V.Kagramanian@iaea.org

**Alexander Soloviev**

International Institute

of Earthquake Prediction Theory

and Mathematical Geophysics

Warshavskoye shosse 79, kor.2

113556 Moscow

RUSSIA

E-mail: soloviev@mitp.ru

**Marcin Paprzycki**

Oklahoma State University

Computer Science Department

700 N.Greenwood

Tulsa

Oklahoma 74106-0700

USA

E-mail: marcin@orca.st.usm.edu

# Part I: Seismic Models and Insurance

# Block models of lithosphere dynamics: approaches and solutions

*Alexander Soloviev and Vyacheslav Maksimov*

## Abstract

The necessity of catastrophe modeling is stipulated both by the essential increase of losses due to recent natural and man-caused hazards and by the lack of reliable real observation data. Earthquakes are considered as an example of unpredictable catastrophic events of a great destructive force. A brief overview of different approaches to mathematical modeling of lithosphere dynamics is presented. Block models are described in details.

# 1. Introduction

The vulnerability of the human civilization to natural dangers is critically growing due to proliferation of high-risk objects, clustering of population, and destabilization of large cities and industrial regions. It is forecasted that the more frequent and larger damages of man-caused and ecological catastrophes can destroy the existing insurance system [3]. This makes the problem of estimation of risks of natural catastrophes to be very important. For the last third of the 20th century the international investigations became more active in the field of development new conceptions concerning the risks of natural catastrophes. In particular, a number of international programs and projects are realized (including the project of the International Institute of Applied Systems Analysis for management of global safety). However, the decisions of the most important global problems of safety are connected with serious difficulties. It is caused among other reasons by vagueness and incompleteness of information and schematic of mathematical apparatus of analysis and forecast [9].

Earthquakes represent typical local catastrophic natural events of a great destructive force. Today a single earthquake may take up to a million lives; cause material damage up to US$1,000,000,000,000, with chain reaction expanding to worldwide economic depression; trigger major ecological catastrophe (e.g. several Chernobyl-type calamities at once); paralyze national defense. In many developing countries the damage from earthquakes consumes all the increase in the GDP. Critically vulnerable became the low seismicity regions, e.g. European and Indian platforms, Eastern US etc.

Seismic risk is a measure of possible damage from earthquakes. Estimation of seismic risk have to facilitate the choice of a wide variety of seismic safety measures, ranging from building codes and insurance to establishment of rescue-and-relief resources. Different representations of seismic risk are required for the choice of different safety measures. Most of the practical problems require to estimate seismic risk for a territory as a whole, and within this territory - separately for the objects of each type: areas, lifelines, sites of the vulnerable constructions etc. The choice of the territory and of the objects is determined by jurisdiction and responsibility of the decision-maker who is using the estimation.

Each concrete representation of seismic risk has to be derived directly from the primary models: of earthquake occurrence; of strong motion caused by a single earthquake; of territorial distribution of population, property, and vulnerable objects; and of the damage caused by an episode of strong motion.

In this study we focus attention on models of earthquake occurrence. Earthquakes as some of other dangers are governed by non-linear systems, which are hierarchical and have intermediate number of degrees of freedom. So far earthquakes are uncontrolled and unpredictable with a sufficient accuracy. The theoretical estimation of statistical parameters of an earthquake flow is a very difficult problem due to absence of an adequate theoretical base. Study of seismicity with the statistical and phenomenological analysis of the real earthquake catalogues has the disadvantage that the instrumental observation data cover, in general, a time interval which is very short, in comparison with the duration of tectonic processes responsible of the seismic activity. Therefore the patterns of the earthquake occurrence identifiable in a real catalogue may be only apparent and may not repeat in the future. The historical data on seismicity are usually incomplete and do not cover uniformly a region under consideration.

We try to overcome these difficulties by means of numerical modeling the seismic process. The synthetic earthquake catalogue obtained by numerical modeling may cover very long time interval that allows us to acquire a more reliable estimation of the parameters of an earthquake flow. In problems of risk estimation the numerical modeling acts as a generator of possible scenarios of catastrophe occurrence.

The paper has the following structure. First, we give a brief overview of mathematical models of lithosphere dynamics. Then we describe block models in detail. The conclusive section of the

paper is devoted to discussion on necessity of numerical parallel algorithms for solving the problem of modeling dynamics of a real system of tectonic plates.

## 2.    Different approaches to modeling seismic processes

The seismic observations show that features of a seismic flow are different for different active regions. It is reasonable to suggest that this difference is due among other factors to contrasts in the tectonic structure of the regions and in main tectonic movements determining the lithosphere dynamics in the regions. The laboratory studies show specifically that this difference is controlled mainly by the rate of fracturing and heterogeneity of the medium and also by the type of predominant tectonic movements [13]. If a single factor is considered it is difficult to detect its impact on features of a seismic flow by using real seismic observations because the seismic flow is impacted by an assemblage of factors some of which could be larger than one under consideration. It is difficult if not impossible to single out the impact of a single factor by analysis of is real seismic observations. This can be overcome by numerical modeling of the processes generating seismicity and studying synthetic earthquake catalogs obtained (see [2, 11, 14]). One more reason to use the models is due to the fact that the study of seismicity with the statistical and phenomenological analysis of the real earthquake catalogs has the disadvantage that the reliable data cover, in general, a time interval of about one hundred years or even less. This time interval is very short, in comparison with the duration of tectonic processes responsible of the seismic activity, therefore the patterns of the earthquake occurrence identifiable in a real catalog may be only apparent and may not repeat in the future, thus excluding any statistical tests. On the other hand, the synthetic catalog obtained by numerical modeling of the seismic process may cover very long time interval that allows us to acquire a more reliable estimation of the parameters of seismic flow.

The following are among the principal features of the lithosphere that should be incorporated into a model for it to be regarded as adequate: interaction of the processes of different physical origin, spatial and temporal scales, hierarchical block or possibly «fractal» structure, and self-similarity in space, time, and energy. The traditional approach to modeling is based on one specific tectonic fault and, often, one strong earthquake in order to reproduce certain seismic phenomena (relevant to this specific earthquake). In contrast, the class of the slider-block and cellular automata models treats the seismic process in the most abstract way, in order to reproduce general universal properties of seismicity, first of all, the Gutenberg–Richter frequency of occurrence law, migration of events, sequence of aftershocks, seismic cycle and so on [7]. The specific and general approaches have their respective advantages and disadvantages. The first approach, which takes into account detailed information on the local geotectonic environment, usually misses universal properties of a series of events in a system of interacting faults. The second approach may be treated as a zero-order approximation to reality. However, the importance of this approach to the earthquake prediction problem lies in the possibility to establish analogs with problems in other sciences and to elaborate a new language for the description of seismicity patterns.

So, mathematical models of lithosphere dynamics developed according to a general approach are tools for the study of the earthquake preparation process and useful in earthquake prediction studies [4]. An adequate model should indicate the physical basis of premonitory patterns determined empirically before large events. Note one more time that the available data often do not constrain the statistical significance of the premonitory patterns. The model can be used also to suggest new premonitory patterns that might exist in real catalogs. Although there is no adequate theory of the seismotectonic process, various properties of the lithosphere, such as spatial heterogeneity, hierarchical block structure, different types of non-linear rheology, gravitational and thermodynamic processes, physicochemical and phase transitions, fluid migration and stress corrosion, are probably relevant to the properties of earthquake sequences.

The qualitative stability of these properties in different seismic regions suggests that the lithosphere can be modelled as a large dissipative system that does not essentially depend on the particular details of the specific processes active in a geological system. For the detailed review of the most important directions of modeling seismic processes, see [5]. Here we dwell on the model where the interaction of tectonic faults is taken into account.


## 3. Detailed description of the block models

The block model of lithosphere dynamics exploits the hierarchical block structure of the lithosphere proposed in [1]. The basic principles of the model are developed, for example, in [4]. According to this model, a seismic region is modeled by a system of absolutely rigid blocks of the lithosphere, which are separated by comparatively thin, weak, less consolidated fault zones, such as lineaments and tectonic faults. In the seismotectonic process all deformations and most earthquakes occur in such fault zones. Relative displacements of all blocks are supposed to be infinitely small with respect to their geometric size. The blocks interact between themselves and with the underlying medium. The system of blocks moves as a consequence of prescribed motion of the boundary blocks and of the underlying medium.

In the model the strains are accumulated in fault zones. This reflects strain accumulation due to deformations of plate boundaries. Of course, considerable simplifications are made in the model, but they are necessary to understand the dependence of earthquake flow on main tectonic movements in a region and its lithosphere structure. This assumption is justified by the fact that for the lithosphere the effective elastic modules in the fault zones are significantly smaller than those within the blocks. The blocks are in viscous-elastic interaction with the underlying medium. The corresponding stresses depend on the value of relative displacement. This dependence is assumed to be linear elastic. The motion of the medium underlying different blocks may be different. Block motion is defined so that the system is in a quasi-static state of equilibrium. The interaction of the blocks along fault zones is viscous-elastic too ("normal state") so far as the ratio of the stress to the pressure remains below a certain strength level. When the critical level is exceeded in some part of a fault zone, a stress-drop ("failure") occurs (in accordance with the dry friction model), possibly causing failure in other parts of the fault zones. These failures produce earthquakes. Immediately after the earthquake and for some time after, the affected parts of the fault zones are in a creep state. This state differs from the normal one because of a faster growth of inelastic displacements, lasting until the ratio of the stress to the pressure falls below some other level. The process of numerical simulation produces a synthetic earthquake catalog as a result.

On the base of idea outlined above a family of block models taking into account real geometry of tectonic regions was developed. The key point for further modifications is so-called two-dimensional plane model the detailed description of which is given below. The paper [12] is devoted to investigation of three-dimensional block movements. In [10] the model is transferred into the sphere in order to simulate global tectonic plate dynamics. To reproduce in a model space-temporal clustering of events, the influence of fluids migrating along tectonic faults was taken into account [16]. The main principles of block models will be described on the example of two-dimensional model as the one, which is more studied than others [6, 8].


## 3.1 Block structure geometry

A layer with thickness $H$ limited by two horizontal planes is considered (Fig. 1), and a block structure is defined as a limited and simply connected part of this layer. Each lateral boundary of the structure is defined by portions of the parts of planes intersecting the layer.

Fig. 1. Block structure: elements and notions

The subdivision of the structure into blocks is performed by planes intersecting the layer. The parts of these planes, which are inside the block structure and its lateral faces, are called "fault zones". The geometry of the block structure is defined by the lines of intersection between the fault zones and the upper plane limiting the layer (these lines are called "faults"), and by the angles of dip of each fault zone. Three or more faults cannot have a common point on the upper plane, and a common point of two faults is called "vertex". The direction is specified for each fault and the angle of dip of the fault zone is measured on the left of the fault. The positions of a vertex on the upper and the lower plane, limiting the layer, are connected by a segment ("rib") of the line of intersection of the corresponding fault zones. The part of a fault zone between two ribs corresponding to successive vertices on the fault is called "segment". The shape of the segment is a trapezium. The common parts of the block with the upper and lower planes are polygons, and the common part of the block with the lower plane is called "bottom". It is assumed that the block structure is bordered by a confining medium, whose motion is prescribed on its continuous parts comprised between two ribs of the block structure boundary. These parts of the confining medium are called "boundary blocks".

## 3.2    Block movement

The blocks are assumed to be rigid and all their relative displacements take place along the bounding fault zones. The interaction of the blocks with the underlying medium takes place along the lower plane, any kind of slip being possible. The movements of the boundaries of the block structure (the boundary blocks) and the medium underlying the blocks are assumed to be an external force on the structure. The rates of these movements are considered to be horizontal and known.

Non-dimensional time is used in the model, therefore all quantities that contain time in their dimensions are referred to one unit of the non-dimensional time, and their dimensions do not contain time. For example, in the model, velocities are measured in units of length and the velocity of 5 cm means 5 cm for one unit of the non-dimensional time. When interpreting the results a realistic value is given to one unit of the non-dimensional time. For example if one unit of the non-dimensional time is one year then the velocity of 5 cm, specified for the model, means 5 cm/year. At each time the displacements of the blocks are defined so that the structure is in a quasistatic equilibrium, and all displacements are supposed to be infinitely small, compared with the block size. Therefore the geometry of the block structure does not change during the simulation and the structure does not move as a whole.

## 3.3   Interaction between the blocks and the underlying medium

The elastic force, which is due to the relative displacement of the block and the underlying medium, at some point of the block bottom, is assumed to be proportional to the difference between the total relative displacement vector and the vector of slippage (inelastic displacement) at the point. The elastic force per unit area $\mathbf{f}^u = (f_x^u, f_y^u)$ applied to the point with co-ordinates $(X,Y)$, at some time $t$, is defined by

$$f_x^u = K_u(x - x_u - (Y - Y_c)(\varphi - \varphi_u) - x_a),$$

$$f_y^u = K_u(y - y_u + (X - X_c)(\varphi - \varphi_u) - y_a), \tag{1}$$

where $X_c$, $Y_c$ are the co-ordinates of the geometrical center of the block bottom; $(x_u, y_u)$ and $\varphi_u$ are the translation vector and the angle of rotation (following the general convention, the positive direction of rotation is anticlockwise), around the geometrical center of the block bottom, for the underlying medium at time $t$; $(x,y)$ and $\varphi$ are the translation vector of the block and the angle of its rotation around the geometrical center of its bottom at time $t$; $(x_a, y_a)$ is the inelastic displacement vector at the point $(X,Y)$ at time $t$.

The evolution of the inelastic displacement at the point $(X,Y)$ is described by the equations

$$\frac{dx_a}{dt} = W_u f_x^u, \qquad \frac{dy_a}{dt} = W_u f_y^u. \tag{2}$$

The coefficients $K_u$ and $W_u$ in (1) and (2) may be different for different blocks.

## 3.4   Interaction between the blocks along the fault zones

At the time $t$, in some point $(X,Y)$ of the fault zone separating the blocks numbered $i$ and $j$ (the block numbered $i$ is on the left and that numbered $j$ is on the right of the fault) the components $\Delta x$, $\Delta y$ of the relative displacement of the blocks are defined by

$$\Delta x = x_i - x_j - (Y - Y_c^i)\varphi_i + (Y - Y_c^j)\varphi_j,$$

$$\Delta y = y_i - y_j + (X - X_c^i)\varphi_i - (X - X_c^j)\varphi_j. \tag{3}$$

where $X_c^i$, $Y_c^i$, $X_c^j$, $Y_c^j$ are the co-ordinates of the geometrical centers of the block bottoms, $(x_i, y_i)$, and $(x_j, y_j)$ are the translation vectors of the blocks, and $\varphi_i$, $\varphi_j$ are the angles of rotation of the blocks around the geometrical centers of their bottoms, at time $t$. In accordance with the

assumption that the relative block displacements take place only along the fault zones, the displacements along the fault zone are connected with the horizontal relative displacement by

$$\Delta_t = e_x \Delta x + e_y \Delta y,$$

$$\Delta_l = \Delta_n / \cos\alpha, \quad \text{where} \quad \Delta_n = e_x \Delta y - e_y \Delta x. \tag{4}$$

That is the displacements along the fault zone are projected on the horizontal plane. Here $\Delta_t$, $\Delta_l$ are the displacements along the fault zone parallel ($\Delta_t$) and normal ($\Delta_l$) to the fault line on the upper plane, ($e_x$, $e_y$) is the unit vector along the fault line on the upper plane, $\alpha$ is the dip angle of the fault zone, and $\Delta_n$ is the horizontal displacement, normal to the fault line on the upper plane. The elastic force per unit area $\mathbf{f} = (f_t, f_l)$ acting along the fault zone at the point $(X,Y)$ is defined by

$$f_t = K(\Delta_t - \delta_t),$$

$$f_l = K(\Delta_l - \delta_l). \tag{5}$$

Here $\delta_t$, $\delta_l$ are inelastic displacements along the fault zone at the point $(X,Y)$ at time $t$, parallel ($\delta_t$) and normal ($\delta_l$) to the fault line on the upper plane. The evolution of the inelastic displacement at the point $(X,Y)$ is described by the equations

$$\frac{d\delta_t}{dt} = Wf_t, \quad \frac{d\delta_l}{dt} = Wf_l. \tag{6}$$

The coefficients $K$ and $W$ in (5) and (6) may be different for different faults. The coefficient $K$ can be considered as the shear modulus of the fault zone.
In addition to the elastic force, there is the reaction force which is normal to the fault zone; the work done by this force is zero, because all relative movements are tangent to the fault zone. The elastic energy per unit area at the point $(X,Y)$ is equal to

$$e = (f_t(\Delta_t - \delta_t) + f_l(\Delta_l - \delta_l))/2. \tag{7}$$

From (4) and (7) the horizontal component of the elastic force per unit area, normal to the fault line on the upper plane, $f_n$ can be written as:

$$f_n = \frac{\partial e}{\partial \Delta_n} = \frac{f_l}{\cos\alpha}. \tag{8}$$

It follows from (8) that the total force acting at the point of the fault zone is horizontal if there is the reaction force, which is normal to the fault zone. The reaction force per unit area is equal to

$$p_0 = f_l \text{tg}\alpha. \tag{9}$$

The reaction force (9) is introduced and therefore there are not vertical components of forces acting on the blocks and there are not vertical displacements of blocks.
Formulas (3) are valid for the boundary faults too. In this case one of the blocks separated by the fault is the boundary block. The movement of these blocks is described by their translation and rotation around the origin of co-ordinates. Therefore the co-ordinates of the geometrical center of the block bottom in (3) are zero for the boundary block. For example, if the block numbered $j$ is a boundary block, then $X_c^j = Y_c^j = 0$ in (3).

## 3.5  Equilibrium equations

The components of the translation vectors of the blocks and the angles of their rotation around the geometrical centers of the bottoms are found from the condition that the total force and the total moment of forces acting on each block are equal to zero. This is the condition of quasi-static equilibrium of the system and at the same time the condition of minimum energy. The forces arising from the specified movements of the underlying medium and of the boundaries of the block structure are considered only in the equilibrium equations. In fact it is assumed that the action of all other forces (gravity, etc.) on the block structure is balanced and does not cause displacements of the blocks.

In accordance with formulas (1), (3)-(5), (8), and (9) the dependence of the forces, acting on the blocks, on the translation vectors of the blocks and the angles of their rotations is linear. Therefore the system of equations which describes the equilibrium is linear one and has the following form

$$A\mathbf{z} = \mathbf{b} \tag{10}$$

where the components of the unknown vector $\mathbf{z} = (z_1, z_2, ..., z_{3n})$ are the components of the translation vectors of the blocks and the angles of their rotation around the geometrical centers of the bottoms ($n$ is the number of blocks), i.e. $z_{3m-2} = x_m$, $z_{3m-1} = y_m$, $z_{3m} = \varphi_m$ ($m$ is the number of the block, $m = 1, 2, ..., n$).

The matrix $A$ does not depend on time and its elements are defined from formulas (1), (3-5), (8), and (9). The moment of the forces acting on a block is calculated relative to the geometrical center of its bottom. The expressions for the elements of the matrix $A$ contain integrals over the surfaces of the fault segments and of the block bottoms. Each integral is replaced by a finite sum, in accordance with the space discretization described in the next section. The components of the vector $\mathbf{b}$ are defined from formulas (1), (3-5), (8), and (9) as well. They depend on time, explicitly, because of the movements of the underlying medium and of the block structure boundaries and, implicitly, because of the inelastic displacements.

## 3.6  Discretization

Time discretization is performed by introducing a time step $\Delta t$. The state of the block structure is considered at discrete values of time $t_i = t_0 + i\Delta t$ ($i = 1, 2, ...$), where $t_0$ is the initial time. The transition from the state at $t_i$ to the state at $t_{i+1}$ is made as follows: (i) new values of the inelastic displacements $x_a$, $y_a$, $\delta_t$, $\delta_l$ are calculated from equations (2) and (6); (ii) the translation vectors and the rotation angles at $t_{i+1}$ are calculated for the boundary blocks and the underlying medium; (iii) the components of $\mathbf{b}$ in equations (10) are calculated, and these equations are used to define the translation vectors and the angles of rotation for the blocks. Since the elements of $A$ in (10) are not functions of time, the matrix $A$ and the associated inverse matrix can be calculated only once, at the beginning of the calculation. Formulas (1-9) describe the forces, the relative displacements, and the inelastic displacements at points of the fault segments and of the block bottoms. Therefore the discretization of these surfaces (partition into «cells») is required for the numerical simulation. It is made according to the special rule, and the co-ordinates $X$, $Y$ and the corresponding inelastic displacements are supposed to be the same for all the points of a cell.

## 3.7  Earthquake and creep

Let us introduce the quantity

$$\kappa = \frac{|\mathbf{f}|}{P - p_0} \qquad (11)$$

where $\mathbf{f} = (f_t, f_l)$ is the vector of the elastic force per unit area given by (5), $P$ is assumed equal for all the faults and can be interpreted as the difference between the lithostatic and the hydrostatic pressure, $p_0$, given by (9), is the reaction force per unit area. For each fault the following three values of $\kappa$ are considered $B > H_f \geq H_s$.

Let us assume that the initial conditions for the numerical simulation of block structure dynamics satisfy the inequality $\kappa < B$ for all the cells of the fault segments. If, at some time $t_i$, the value of $\kappa$ in any cell of a fault segment reaches the level $B$, a failure ("earthquake") occurs. The failure is meant as slippage during which the inelastic displacements $\delta_t$, $\delta_l$ in the cell change abruptly to reduce the value of $\kappa$ to the level $H_f$. Thus, the earthquakes occur in accordance with the dry friction model. The new values of the inelastic displacements in the cell are calculated from

$$\delta_t^e = \delta_t + \gamma f_t, \quad \delta_l^e = \delta_l + \gamma f_l \qquad (12)$$

where $\delta_t$, $\delta_l$, $f_t$, $f_l$ are the inelastic displacements and the components of the elastic force vector per unit area just before the failure. The coefficient $\gamma$ is given by

$$\gamma = 1/K - PH_f/(K(|\mathbf{f}| + H_f f_l \mathrm{tg}\alpha)) \qquad (13)$$

It follows from (5), (9), (11-13) that after the calculation of the new values of the inelastic displacements the value of $\kappa$ in the cell is equal to $H_f$. After calculating the new values of the inelastic displacements for all the failed cells, the new components of the vector $\mathbf{b}$ are calculated, and from the system of equations (10) the translation vectors and the angles of rotation for the blocks are found. If for some cell(s) of the fault segments $\kappa > B$, the procedure given above is repeated for this cell (or cells). Otherwise the state of the block structure at the time $t_{i+1}$ is determined as follows: the translation vectors, the rotation angles (at $t_{i+1}$) for the boundary blocks and for the underlying medium, and the components of $\mathbf{b}$ in equations (10) are calculated, and then equations (10) are solved.

Different times could be attributed to the failures occurring on different steps of the procedure: if the procedure consists of $p$ steps the time $t_i + (j - 1)\delta t$ can be attributed to the failures occurring on the $j$th step, and the value of $\delta t$ is selected to satisfy the condition $p\delta t < \Delta t$. The cells of the same fault zone in which failure occurs at the same time form a single earthquake. The parameters of the earthquake are defined as follows: (i) the origin time is $t_i + (j - 1)\delta t$; (ii) the epicentral co-ordinates and the source depth are the weighted sums of the co-ordinates and depths of the cells included in the earthquake (the weight of each cell is given by its square divided by the sum of squares of all the cells included in the earthquake); (iii) the magnitude is calculated from the formula [15]:

$$M = D\mathrm{lg}S + E, \qquad (14)$$

where $D$ and $E$ are constants and $S$ is the sum of the squares of the cells (in km$^2$) included in the earthquake. Immediately after the earthquake, it is assumed that the cells in which a failure has occurred are in the creep state. It means that, for these cells, in equations (6), which describe the evolution of inelastic displacement, the parameter $W_s$ ($W_s > W$) is used instead of $W$, and $W_s$ may be different for different faults. After each earthquake a cell is in the creep state as long as $\kappa > H_s$, while when $\kappa \leq H_s$, the cell returns to the normal state and henceforth the parameter $W$ is used in (6) for this cell.

## 4. Parallel algorithm for numerical simulation

Computational experiments showed that the block models of lithosphere dynamics during performing on sequential computers require considerable expenditures of memory and time of a processor, and it does not allow to simulate dynamics of complicated structures.

However, the approach applied to modeling admits sufficiently effective parallelization of calculations on a multiprocessor machine, and it makes real passing to a system of tectonic plates in the global scale (with the use of real geophysical and seismic data) and to the spherical geometry [10].

On working stations basing on microprocessors Alpha-21164 (533MHz, 256Mb) at IMM UB RAS (Ekaterinburg, Russia) the variant of parallel program was realized by the scheme «master–worker» («processor farm»). The demands of compatibility with different platforms (in the sense of fast transition, ideally, by means of simple recompiling) were made to the program code. For this purpose, the special library MPI («message passing interface») was used, and the parallel algorithm was designed in such a way that the unique loading module was formed for all processors. The block-scheme of this algorithm is presented in Fig. 2-4. Let us give necessary explanations.

In the beginning of the work the number of processor the program has loaded to is detected (zero processor becomes the master). After this process, the information on a block structure is red, and auxiliary calculations (before the main cycle) are performed. It is important that a part of calculations performed only by the master (due to finding block and underlying medium displacements according to (10)) requires insignificant time expenditures. At every time step the most time-consumable procedure is calculation of values of forces and inelastic displacements in all cells of space discretization of block bottoms and fault segments. Since these calculations may be performed independently from each other, they are shared between all processors, each of which processes own portion of cells.



Fig. 2. Scheme of parallelization of the block model. Notation: operations carried out only by master are marked by «M», only by workers – by «W».

The exchange of information between processors at every time step is realized according to the following scheme. The master calculates new values of block, boundary block and underlying medium displacements, then necessary parameters are transferred to the workers. Recalculated values of the right-hand part of system (10) are returned to the master, then the next time step is carried out. For processing the situation treated as an earthquake (section 3.7), the scheme is slightly complicated, since in this case the master should ask all the workers until cells of segments in the critical state exist. The time of calculations on each processor is much more than the time of exchange. Therefore rather high useful loading of each processor is achieved.

```
                            START

Distribution of cells of space discretization into portions
          (depending on the number of processors)

Cycle with respect to time: calculation of ξ

              Procedure CALC

Finding elastic forces and stress on segments (5), (6), (8)

    TRANSFER of maximum stress                W
    RECEIVEING                                M

    TRANSFER of the flag «earthquake!» (11)   M
    RECEIVING                                 W

Are there              yes    Processing
earthquakes?                  earthquakes, creep
                              (11)-(13)
      no

Finding elastic forces on block bottoms (1), (2)

    TRANSFER of data on state of cells        W
    RECEIVING, writing into text file         M

Is there end of        no
the cycle?
      yes

    TRANSFER of parameters to continue calculation   W
    RECEIVING, saving in file, stopping workers, output of
    results                                           M

                            END
```

Fig. 3. Procedure *RUN*

For testing the dependence of time of solving the problem on the number of processors and comparing with sequential algorithm, the following values were analyzed: acceleration coefficient $S_r = T_1/T_r$ and effectiveness coefficient $E_r = S_r /r$, where $T_r$ is the time of program performance on multiprocessor computer with $r$ processors, $T_1$ is the corresponding time for sequential algorithm. Note that $T_r$ is the sum of pure time of calculations and expenditures for necessary exchanges. It is appeared that $S_r$ is slightly less than $r$, consequently, $E_r$ is close to 1, and the parallelization effectiveness is rather high and it insignificantly decreases with increasing the number of processors in action (in correspondence with the parallelization scheme).



Fig. 4. Procedure *CALC*

The scheme described in this section was applied to simulation of dynamics of different block structures: both model and approximations of real regions. However, presentation of results of modeling is out of the scope of this paper (see, for example, [10]).

## References

1. Alekseevskaya, M.A., Gabrielov, A.M., Gvishiani, A.D., Gelfand, I.M. and E.Ya.Ranzman, 1977, Formal morphostructural zoning of mountain territories, *J. Geophys. Res.,* **43,** pp.227–233.

2. Burridge, R., and Knopoff, L., 1967, Model and theoretical seismicity, *Bull. Seismol. Soc. Amer.,* **57**, pp.341–371.

3. Ermoliev, Yu.M., Ermolieva, T.Y., MacDonald, G.J. et al., 2000, A system approach to management of catastrophic risks, *Eur. J. Oper. Res.* No. 122, pp.452–460.

4. Gabrielov, A.M., Levshina, T.A., and Rotwain, I.M., 1990, Block model of earthquake sequence, *Phys. Earth and Planet. Inter., **61,** pp.18–28.

5. Gabrielov, A.M., 1993, Modeling of seismicity, Second Workshop on Non-Linear Dynamics and Earthquake Prediction, 22 November – 10 December, 1993, Trieste, Italy. Preprint, 22 p.

6. Gorshkov, A., Keilis-Borok, V., Rotwain, I., Soloviev, A., and Vorobieva, I., 1997, On dynamics of seismicity simulated by the models of blocks-and-faults systems, *Annali di Geofisica,* **XL,** 5: pp.1217–1232.

7. Kagan,Y., and Knopoff, L., 1978, Statistical study of the occurrence of shallow earthquakes, *Geophys. J. R. Astron. Soc.*, **55,** pp.67–86.

8. Keilis-Borok, V.I., Rotwain, I.M., and Soloviev, A.A, 1997, Numerical modeling of block structure dynamics: dependence of a synthetic earthquake flow on the structure separateness and boundary movements, *Journal of Seismology*, **1,** 2: pp.151–160.

9. Marchuk, G.I., and Kondratiev, K.Ya., 1992, Problems of global ecology. M.: Nauka, 264 p.

10. Melnikova, L.A., Rozenberg, V.L., Sobolev, P.O., and Soloviev, A.A., 2000, Numerical simulation of dynamics of a system of tectonic plates: spherical block model, *Comp. Seismology,* Iss.31, pp.138–153.

11. Newman, W.I., Turcotte, D.L., and Gabrielov, A.M., 1995, Log-periodic behaviour of a hierarchical failure model with application to precursory seismic activation, *Phys. Rev. E., **52,** pp.4827–4835.

12. Rozenberg, V., and Soloviev, A., 1997, Considering 3D Movements of Blocks in the Model of Block Structure Dynamics. Fourth Workshop on Non-Linear Dynamics and Earthquake Prediction, 6–24 October, 1997, Italy. Preprint, 26 p.

13. Sherman, S.I., Borniakov, S.A., and Buddo, V.Yu., 1983, Areas of Dynamic Effects of Faults. Novosibirsk: Nauka (in Russian).

14. Turcotte,D.L., 1997, Fractals and Chaos in Geology and Geophysics. 2nd Ed., Cambridge University Press.

15. Utsu, T., and Seki, A., 1954. A relation between the area of aftershock region and the energy of main shock, *J. Seism. Soc. Japan, **7,** pp.233–240.

16. Zheligovskii V.A., Podvigina, O.M., and Gabrielov, A.M., Migration of fluids and dynamics of a block-and-fault system, *Comp. Seismology,* Iss.33 (in press).

# Risk insurance:  Generation of scenarios[1]

*Marina Blizorukova, Tatiana Ermolieva,*
*Valerii Rozenberg (in cooperation with Valerii Samosyuk)*

## Abstract

Problems of insurance against natural hazards and related modeling tools are discussed. Methods of generation of possible scenarios of earthquakes are in the focus. It is conjectured that the generated scenarios may act as inputs to optimization algorithms aimed at indicating possible improvements in the structure of the regional insurance networks. A review on relevant modeling approaches is given and their applicability for different classes of data discussed.

---

# 1. Introduction: the role of catastrophe modeling

The tendencies in the socio-economic development and environmental global changes, which have become a dominant feature of the recent decades, have led to a dramatic and rapid increase of losses due to natural and anthropogenic catastrophic events. Within the last three decades the direct catastrophe damages only from natural disasters have increased nine-fold [4]. Catastrophes destroy communication systems, electricity supply and irrigation, they affect consumption, savings and investments. It should be noted that low-income countries with transition economics are especially sensitive to such losses. One of the main reasons for the increase of catastrophe damages is the ignorance of risks leading to the clustering of people and capital in the risk prone areas as well as the creation of new risk prone areas. It is estimated [13] that within the next fifty years more than a third of the world population will live in seismically and volcanically active zones. This alarming human-induced tendency calls for new risk-based advanced computational approaches to economic and insurance developments. In this paper we focus on such an important aspect of large-scale problems of decision making on ex-ante risk reduction measures and loss spreading mechanisms as insurance contracts against natural hazards which begin to play a significant role in managing catastrophe losses (see, for example, [1-5]). These decisions are evaluated using so-called catastrophe modeling. There is a number of methodological challenges involved in catastrophic risk management [4]. Here, we outline its characteristic features such as endogenous risks, mutually dependent losses, the lack of information, the need for long-term perspectives and geographically explicit models, and others.

## 1.1 Complex interdependencies

Catastrophes produce severe losses characterized by mutual dependence in space and time. The multivariate distribution of these losses is, in general, intractable analytically. It depends on the clusterization of values in the region and on the patterns of catastrophes. Besides, it may dramatically depend on policy variables. For example, a dam fundamentally modifies flood conditions downstream and along the site. This creates favorite conditions for insurance and new land-use transformations. On the other hand, a failure of the dam may lead to rare but more devastating losses in the protected area. Such interdependencies of decisions and risks restrict the straightforward "one-by-one" evaluations of feasible options. The so-called "if-then" analysis runs quickly into an extremely high number of alternatives. Thus, with only 10 feasible decisions, say 10%, 20%,…, 100% of the insurance coverage for a particular site, and 10 possible heights of the dam, the number of possible "if-then" combinations is $10^{10}$. At one second per evaluation, more than 90 years are required to carry out the computations. The main idea in dealing with this problem is to avoid exact evaluations of all possible alternatives and concentrate attention on the most promising directions. From a formal point of view this is equivalent to the design of special search techniques (in the space of decision variables), making use of random simulations of catastrophes. This is a task of stochastic optimization [6]. Certain of these search procedures can also be viewed as adaptive scenario analysis, or adaptive Monte Carlo optimization [2-3]. They generate feedback to policy variables after each simulation and automatically drive them towards desirable combinations without going into exhausting "if-then" analyses.

## 1.2 Rare events

The principal problem with the management of rare catastrophic risks is the lack of historical data on losses at any particular location, although rich data may exist on an aggregate regional

level. Historical data are relevant to old policies and may have very limited value for new policies. Models have to play a key role for generating data and designing new policies.

Catastrophes may be of quite different nature from episode to episode, exhibiting a wide spectrum of impacts on public health, the environment and the economy. Each of these episodes seems to be improbable and may be simply ignored in the so-called "practical approaches" or "scenario thinking". This may lead to rather frequent "improbable" catastrophes: although each of $N$ scenarios episodes) has a negligible probability $p$, the probability of one of them increases exponentially in $N$ as $1-(1-p)^N = 1-\exp(N \ln(1-p))$. In other words, the integrated analysis of all possible, although rare, scenarios is essential.

## 1.3   Long-term perspectives

The proper assessment and management of rare risks requires also long-term perspectives. The occurrence of a catastrophe within a small interval $\Delta t$ is often evaluated by a negligible probability $\lambda \Delta t$, but the probability of a catastrophe in an interval [0,T] increases as $1-(1-\lambda \Delta t)^{T/\Delta t} \approx 1-e^{-\lambda T}$. Purely adaptive "learning-by-doing" or "learning-by-catastrophe" approaches may be extremely expensive. The year-by-year adjustments of economic developments with so-called anualization of catastrophes may be very misleading. In this case a 50-years catastrophe of an airplane is reduced, in fact, to the sum of annual crashes of its parts, say, wheels in the first year, a wing in the second, and so on.

## 1.4   Spatial aspects

Catastrophes have different spatial patterns and quite differently affect locations. For example, the location of properties or structures with regard to the center of an earthquake is an extremely important piece of information. Together with the regional geology and the soil conditions the location influences the degree of shaking, and, hence, damage incurred at the location. The deforestation at a particular location modifies the flood conditions only downstream and affects the insurance claims only from specific locations. In other words, management of complex interdependencies among catastrophic risks, losses and decisions is possible only within a geographically explicit framework.

## 1.5   Robust management strategies

Uncertainty is associated with every facet of catastrophic risk assessment. The exact evaluation of all complex interdependencies is impossible and thus risk assessment will yield poor estimates. In this situation the most important task seems to be the design of robust management strategies. Although the assessment is not exact, the preference structure among different decisions may be rather stable to errors. This is similar to the situation with two parcels: to find out their weights is a much more difficult task than to determine the heavier parcel. This simple observation, in fact, is the basic idea of stochastic optimization approaches proposed in [1-4], namely, the evaluation of the optimal decision is achieved without exact evaluation of all possible alternatives.

## 1.6    Multiagent aspects

The high consequences of catastrophes call for the cooperation of various agents such as governments, insurers, investors, and individuals. This often leads to multi-objective stochastic optimization problems and game-theoretical models with stochastic uncertainties [6].

For all these reasons models become essential for catastrophic risks management. The occurrence of various episodes (scenarios) and dependent losses in the region can be simulated on a computer in the same way as the episode may happen in reality [14]. The stochastic optimization techniques can utilize this information for designing robust management strategies.

## 2.    Modeling earthquakes as a decision making tool in insurance

In this paper we focus on some aspects of catastrophe modeling and integrated management in the case of earthquakes which represent typical local natural catastrophic events of a great destructive force. So far these phenomena are not well understood from a physical standpoint, they are uncontrolled and unpredictable with a sufficient accuracy. Reliable statistical analysis of earthquakes is rather difficult since existing observation data cover only short time intervals. Nevertheless it became clear that strategies for insuring property against such events can be based on catastrophe modeling to compensate the lack of real information on possible damages and the absence of analytical representation. All models should contain, rougly speaking, three modules: the seismic hazard module (earthquake scenario generator), the vulnerability module, and the financial module (their characteristic features were mentioned in items 1.1-1.5 of the previous section).

The earthquake scenario generator simulates actual earthquake shaking. This module often comprises other physical phenomena associated with an earthquake including subsequent fires, landslides. The movement of seismic waves through the soil is modeled by attenuation equations. Seismic effects at a site depend on earthquake magnitude, intensity, deep, distance from the epicenter, and site characteristics, such as regional geology and soil types. As a rule, the earthquake scenario generator should produce a sequence of events (earthquake catalog). The analysis of synthetic events on time intervals of arbitrary lengths may reveal correlations in the occurrence of events and detect phenomena preceding strong catastrophes.

The vulnerability module relates seismic shaking to structural and property damage. It determines the extent of damages to buildings and content at a site.

The financial module assigns a cost to these damages and calculates the maximum potential and/or expected losses for either individual sites and regions. It calculates losses due to structural damage, damage to property and content, and often business interuption. This includes data on building locations, type and contents. The estimates are presented either in percentage of the total value or as a monetary value.

The histograms of aggregate losses for a single location, a particular catastrophe zone or a country can be derived from catastrophe modeling. But it has only marginal benefits when it is used in a traditional manner for obtaining estimates of aggregate losses. First of all, this type of modeling is a decision making tool, but the decision variables are not explicitly incorporated in the existing catastrophe models. Following [13], we admit that the currently existing form of catastrophe modeling can only be a necessary subset of more extensive models used to optimize portfolios of risks in an integrated manner.

So, modeling sequences of earthquakes can be treated as a basis for efficient planning mitigation meausures and for insurance strategies in seismic regions. The general scheme of construction of optimal insurance network by means of a scenario-based approach is presented in Fig.1. But the essence of this approach as well as the statement of optimization problem are

out of the scope of this paper (refer to [1, 3] for details). Here we focus our attention only on the problem of generation of earthquake scenarios.



Fig.1. Flowchart of the approach to construction of optimal insurance network.

Depending on the character of data available for a region under investigation, one of the two approaches to modeling earthquakes can be implemented. The first, stochastic, approach is based on a special statistical analysis of geophysical and seismic data. The second one uses results of simulation by means of adequate mechanical models. In both cases, the output is a sequence of possible earthquakes in a given region, which, being combined with the information on the vulnerability of buildings and on the costs of the regional property, serves as one of the input parameters for an optimization model responsible for the design of robust insurance decisions (see, for example, [1- 3]), as it is illustrated in Fig.1.

The goal of this paper is to consider possibilities to generate earthquake scenarios for the concrete seismoactive region (we take the South American region as an example). It turns out that the data available are not sufficient to apply the first approach to the region successfully, so we give only its outline. We concentrate on the second approach although it is under permanent development too.

## 3. Stochastic approach

This approach is based on statistical models analyzing real geophysical and seismic data available for some particular region. The following input data are used to generate an earthquake scenario: a map of seismic activity zones, a map of maximum observed macroseismic intensities, a map of the geotectonic structure of the region, the Gutenberg-Richter relation for the region. Using the map of seismic activity zones, one can estimate a probability of the occurrence of an earthquake at a given point for a given expectation time. This estimation provides a basis for Monte-Carlo simulations of the occurrence of epicenters. It is not a reliable way to find the magnitude and intensity of a strongest possible earthquake at a fixed point using available seismic observations since extreme events occur too rarely. To find a possible intensity, one can use the method based on Gumbel's extreme values distribution, for example, of the form [7]:

$$F(I) = 1 - \exp[\exp(\alpha I_S + \beta) - \exp(\alpha I + \beta)].$$

18

Here, $I$ is the intensity; parameters $I_S$, $\alpha$, and $\beta$ are derived individually for every particular region.

A possible magnitude of a model earthquake may be selected with the use of the available catalog of events. If the number of records in the catalog is not too large, the maximum magnitude observed in the past at a particular point is taken as the magnitude of the corresponding model earthquake. One can also use the statistical method described in [12]. This method is based on the assumption that the following Gutenberg-Richter frequency-magnitude relation holds true:

$$\log N(M) = a^* - b^* M, \quad M_0 \leq M \leq M_{\max}$$

where $N(M)$ is the average number of events with a magnitude no less than $M$ per a unit of time (e.g., one year) in the given region; $[M_0, M_{\max}]$ is the interval for the magnitude values; $a^*$, $b^*$ are parameters.

In accordance with the formula above, the random variable $M$ is described by the following distribution function:

$$F(M) = \frac{10^{-bM_0} - 10^{-bM}}{10^{-bM_0} - 10^{-bM_{\max}}}.$$

It should be noted that the geotectonic structure of the region determines a specific type of the distribution functions given above. They are used to create a sequence of possible earthquakes according to the standard Monte-Carlo procedure.


## 4.   Mechanical approach

The usage of the approach outlined above is not applicable in the case when data available for some region cover a relatively short time interval (the lack of reliable information).In this situation another approach to generate earthquake scenarios is suggested. Briefly, it consists in the following. We use block models of lithosphere dynamics [8, 9] to obtain a sequence of synthetic earthquakes for the region under consideration. In these models a seismically active region is considered as a system of absolutely rigid blocks separated by infinitely thin plane faults. The motion of the system of blocks is determined by a prescribed motion of the boundaries and the underlying medium. Displacements of the blocks are determined so that the system remains in a quasistatic equilibrium state. Block interactions along the faults are viscous-elastic while the ratio of the stress to the pressure is below a certain strength level. When the level is exceeded for a part of a fault, a stress-drop (a failure) occurs in accordance with the dry friction model. The failures represent earthquakes. A synthetic earthquake catalog is produced through numerical simulations. Every model event from this catalog is characterized by some origin time, epicentral coordinates and depth, magnitude and intensity. All such events occurred in a specified time interval represent one earthquake scenario. The number of possible earthquake scenarios depends on the ratio of the length of the whole time interval taken for modeling and the length of the interval specified for scenario generation. From the viewpoint of the problem under consideration the important feature of block models is the possibility to simulate earthquake sequences on arbitrary long time intervals, so we can obtain arbitrary number of earthquake scenarios. It is evident that a model applied to a region should be adequate in the sense that it should reproduce main patterns and features which were determined empirically in real seismic flow in this region (the Gutenberg-Richter law, periodicity of strong events, clustering of events and so on). The basic constructions and ideas of block models of

lithosphere dynamics are described in detail, for example, in [8]. The first results of application of the mechanical approach to modeling possible earthquakes in the South American region are presented in [5].


## 5.  Simulation results

The earthquake sequences obtained according to the scheme outlined above are input parameters for an earthquake scenario generator. We started our numerical experiments with modeling the following subsystem of plates by means of a spherical modification of the block model. The structure includes South America, Caribbean, Cocos, and Nazca plates (Fig.2). Other, surrounding, plates (North America, Africa, Antarctica, and Pacific) are treated as boundary blocks whose motions are prescribed [10]. This region was chosen because of the lack of empirical information, which does not allow to generate scenarios using the statistical approach. The structure under consideration has 4 blocks, 33 vertices, 36 faults, and 4 boundary blocks. Dip angles of faults at boundary South America/Nazca equal $50^0$, other faults have dip angles of $90^0$. The time discretization step was 0.01, and space discretization step 3 km. for segments and $1/3^0$ for block bottoms. The largest block's bottom was split into 40000 cells.

The block models of lithosphere dynamics (especially the spherical modification) are quite time and memory consuming on sequential computers which does not allow to model the dynamics of complex structures. However, these models admit a natural parallelization on multi-processor machines. High efficient parallelization was engaged in computational procedures.



Fig.2. Results of simulation of plate motion and spatial distribution of strong earthquakes: the directions of model plate motion (arrows), subduction zones (light shading), spreading zones (dark shading), epicenters of model events (asterisks). Numbers stand for the plates: 1 – Nazca, 2 – South American, 3 – Cocos, 4 – Carribean, 5 – North American, 6 – Pacific, 7 – Africa, 8 – Antarctica. Symbol "+" marks the segment with evidently detecting clustering of model events.


The program employed gave us quantitative characteristics of block displacements, which can be treated as velocities (in cm/yr.), and relative displacements of points belonging to the fault segments separating the blocks; the displacements result from the interactions of the tectonic plates. A comparison of the obtained synthetic data with the real ones showed that the zones of subduction and spreading are simulated properly (see Fig.2). Some features of the obtained

synthetic catalog of earthquakes (frequency-magnitude dependencies, spatial distribution of epicenters, clustering phenomenon, and others) inherit the real ones. The Gutenberg-Richter frequence-of-occurrence curves for the synthetic and real catalogs have similar slopes at their most informative parts (see Fig.3).



Fig.3. The frequence-of-occurrence curves constructed for the real catalog (solid line) and synthetic catalog (dashed line); $N$ is the accumulated number of earthquakes, $M$ is the magnitude).

Clustering of events are seen both on separate segments (see Fig.4) and in the whole structure; as a rule, forshocks, main shocks and aftershocks are indicated in groups. The pattern of seismicity is repeated qualitatively over a certain time interval (depending on the fault). Periods of post-seismic relaxation and stress accumulation are also observed.



Fig.4. The dependence of the magnitude of model earthquakes on time for the segment marked by "+" in Fig.2.

One can observe the phenomenon of earthquake migration along the faults. Moreover, spatial distribution of events shows that there are faults, at which a significant part of synthetic seismicity is concentrated (in Fig.2 such spots are marked). These faults belong to the main

21

seismoactive zones (Nazca/South America and Nazca/Pacific boundaries). These facts show that the model is, qualitatively, adequate and can be used as an earthquake scenario generator. The synthetic catalog consists of more than 13000 events during the time period of 20 units of non-dimensional time. Each event occurred in a given area (or each strong event followed by a series of aftershocks, see Fig.4) can be treated as a separate scenario. An earthquake scenario includes the coordinates of the epicenter, intensity and magnitude of the earthquake. To generate a shaking intensity, one can use the structure of the lineaments characterizing the geological structure of a region and a model of isoseists [11]. The lineaments are usually identified with faults representing fracture zones on the Earth's surface. An isoseist $A_I$ is a domain in which the shaking intensity is no less than $I$. The simplest model of an isoseist domain is an ellipse. The ratio between its axes is determined by the parameters of the region, magnitude and the earthquake generation mechanism. The affected area increases as the magnitude grows. Models of isoseists are constructed on the basis of the estimates of the average radii of isoseists for well studied earthquakes. A relation between the magnitude $M$ and shaken area $Q$ (in sq.km.) is represented in the form the linear regression

$$\log Q_I(M) = d_I + f_I M.$$

Parameters $d$ and $f$ are constants depending on the region. As a rule, the ratio between the axes of isoseists is set 1:1.5. The main axes of the isoseist ellipses go along the lineaments mentioned above. An example of a system of model isoseists is presented in Fig.5. The systems of isoseist domains associated to every model event in the region are important output parameters of the earthquake scenario generator.



Fig.5. Isoseist domains for magnitude $M$=7.0. Intensity (in points of the seismic scale): 9 – color 1, 8 – color 2, 7 – color 3.

Using all above information and applying vulnerability and financial modules with necessary data, estimates of possible losses are formed and transported, as inputs, to an optimization procedure (see Fig.1). The larger set of scenarios is processed in this way, the more effectively the optimization procedure works [1-3].

The methodology outlined above is expected to be applied for other seismoactive regions, first of all, the Vrancea region.

# References

1. Digas, B.V., Ermoliev, Yu.M., and Kryazhimski, A.V., 1998, Guaranteed optimization in insurance of catastrophic risks. IIASA Interim Report, IR-98-082.

2. Ermoliva, T.Y., Ermoliev, Yu.M., Norkin, V.I., 1997, Spatial Stochastic Model for Optimization Capacity of Insurance Networks Under Dependent Catastrophic Risks: Numerical Experiments. IIASA Interim Report, IR-97- 028.

3. Ermoliev, Yu.M., Ermolieva, T.Y., MacDonald, G.J., Norkin, V.I., Amendola, A., 2000, A system approach to management of catastrophic risks. IIASA Research Report, PR-00-08 (Reprinted from European Journal of Operation Research, 2000, 122, pp. 452–460).

4. Ermoliev, Yu.M., Ermolieva, T.Y., MacDonald, G.J., and Norkin, V.I., 2000, Catastrophic risk management and economic growth. Proceedings of the Second EuroConference "Global change and catastrophe risk management: earthquake risks in Europe", IIASA, Laxenburg, Austria, 6–9 July 2000.

5. Ermoliev, Yu.M., Soloviev, A.A., Maksimov, V.I., Rozenberg, V.L., and Digas, B.V., 2000, Insurance against earthquakes: different approaches to scenario generation. Proceedings of International Symposium on Science, Research, and Education, Zielona Gora, Poland, 28–29 September 2000, pp. 51–57.

6. Ermoliev, Yu.M., Wets, R. (eds), 1988, Numerical techniques of stochastic optimization. Computational Mathematics, Springer–Verlag.

7. Freudenthal, A.M., and Gumbel, E.J., 1959, Physical and statistical aspects of fatigue, pp. 117–158.

8. Gabrielov, A.M., Keilis-Borok, V.I., Levshina, T.A., and Shaposhnikov, V.A., 1986, Block model of lithosphere dynamics. *Mathematical Methods in Seismology and Geodynamics, Computational Seismology*, vol. 19, Moscow, Nauka, pp. 168–178 (in Russian).

9. Gorshkov, A.I., Keilis-Borok, V.I., Rotwain, I.M., Soloviev, A.A., and Vorobieva, I.Yu., 1997, On dynamics of seismicity simulated by the models of blocks-and-faults systems. *Annali di Geofisica*, XL, No. 5, pp. 1217–1232.

10. Gripp, A.E., and Gordon,R.G., 1990, Current plate velocities relative to the hotspots incorporating the Nuvel-1 global plate motion model. *Geoph. Res. Let.*, vol. 17, No. 8, pp. 1109–1112.

11. Keilis-Borok, V.I., Molchan, G.M., Gotsadze, O.D., Koridze, A.H., and Kronrod, T.L., 1984, Experience of estimation of seismic risk for inhabited buildings in countryside of Georgia, *Computational Seismology*, vol. 17, Moscow, Nauka, pp.58–67.

12. Pisarenko, V.F., 1996, Statistical estimation of parameters related to earthquakes of maximum possible strength. *Journal of Earthquake Prediction Research*, No.5, pp. 194–201.

13. Rundle, J.B., Turcotte, D., and Klein, R. (eds), 1996, Reduction and protection of natural disasters. Addison–Wesley.

14. Walker, G., 1997, Current developments in catastrophe modeling, in Britton N.R. and Oliver J. (eds) "Financial risks management for natural catastrophes", Brisbane, Griffith University, Australia, pp. 17–35.

# Earthquake risk management via stochastic optimization: a case study for an Italian Region

*Aniello Amendola[1], Yuri Ermoliev, Tatiana Ermolieva*

## Abstract

IIASA's research focuses, in particular, on issues of efficiency and equity for disaster loss mitigation and sharing. Within international co-operation networks, the activities are carried out mainly in the form of case studies. This paper describes how a spatial-dynamic stochastic optimization model that takes into account complexities and interdependencies of catastrophic risks has been customized to explicitly incorporate the geological characteristics of a region, parameters of seismic hazards and the vulnerability of the built environment. In its general form, the model can analyse the interplay between investment in mitigation and risk-sharing measures. In its application, special attention is given to the evaluation of a multipillar earthquake loss spread program involving central government and mandatory insurance. The model is shown to be able to analyse multiple policy options for developing insurance in an equitable and fair manner, and their effects on the insurance premium and reserve funds. To analyze the stability of the system, we use a strong connection between the nondifferentiable stochastic optimization and such risk measures as Value-at-Risk (VaR) and the probability of bankruptcy (insolvency).

---

[1] [1] Work performed during a visiting period from EC-JRC-ISIS, Ispra (VA) Italy

# 1. Introduction

On-going investigations at IIASA concern the modelling of catastrophe events (the evolution of economic losses), risk management and policy analysis. As the victims relate these losses to human culpability (because of both inadequate prevention measures and regulation, and moral hazards or negligence in their implementation), issues of equity and efficiency for preventing and absorbing the losses are becoming pre-eminent. These policy issues cannot be addressed solely with expertise on the physical and economic phenomena, but they require an understanding of the diverse social concerns and complex institutional processes.. As discussed in a recent paper (Linnerooth- Bayer & Amendola 2000/a), different views of fairness and equity stemming from different forms of social organisation suggest that neither the market nor the government will be acceptable as the mechanism for disaster burden sharing. Thus, some form of a public-private partnership may be appropriate (Kunreuther & Roth 1998). The methodological challenge is to develop tools and models for assessing geographically and temporally dependent catastrophic risks so that policies for hazard reduction, loss mitigation and loss-sharing strategies can be framed.

A well-known example of a government acting as a primary insurer is the U.S. National Flood Insurance Program (NFIP), which seeks to provide insurance at actuarially fair premiums combined with incentives on communities and homeowners to take appropriate loss-reducing measures (Pasterick 1998). Given the size of the U.S. and the large number of persons living in flood plains, the program is sufficiently diversified to cover most regional losses with premium payments. In contrast to the NFIP, some government insurance schemes in Europe, e.g., the French national insurance program, cross-subsidise claims. This is because the Constitution (1946, 1958) established the principle of "the solidarity and equality of all French citizens facing the expenses incurred through national calamities" (Gilber & Gouy 1998).

However, even if many governments are pursuing policies to reduce their role in compensating uninsured victims, a comprehensive recent study confirms that the victims and their governments bear the major losses from natural disasters (Linnerooth-Bayer & Quijano 2000), and, world-wide, there is only moderate risk-transfer with insurance. For the cases considered in this study, non-reimbursed losses for all sources ranged between 40% and 60% of the estimated direct losses. In countries with a potential for very large losses, market opportunities appear to be enhanced by public involvement in national insurance systems, and insurance companies can and should play a more active role in designing these systems. An important consideration for national insurance strategies is linking private insurance with mitigation measures to reduce losses. Insurers, however, are reluctant to enter markets that expose them to a risk of bankruptcy. In the U.S., for example, many insurers pulled out of catastrophic risk markets in response to their large losses from natural catastrophes in the last decade (Cummins & Doherty 1996).

To reduce their risk of insolvency, insurers' strategies might be based on modelling tools that account for the complexity implied by the manifold dependencies in the stochastic process of catastrophic events, decisions and losses. To study the problem in its complexity a spatial-dynamic, stochastic optimisation model has been developed at IIASA (Ermoliev et al. 1997, Ermolieva, 1997, Ermoliev et al. 2000). The model is based on Monte Carlo simulations of catastrophic events in the selected regions. The key feature of the model is the search technique allowing for adaptive adjustments of decision variables towards desirable outcomes on the basis of sequential simulations. The model makes use of a connection between nondifferentiable (possibly convex) stochastic optimisation and such measures as Value-at-Risk (VaR) and probability of bankruptcy – the key indicators of the insurance business performance.

In a first application, the model analysed the insurability of risks in the Irkutsk region in Russia, which is exposed to the risks of earthquakes (Amendola et al. 2000). The results demonstrated the adequacy of the model to generate insurance strategies that are robust with respect to dependencies and uncertainties, thus reducing the risk of bankruptcy to the insurers.

This paper describes a second application that has been devoted to exemplify earthquake risk management in Italy, where a law for integrating insurance in the overall risk management process was only proposed in late 1997 (within the Design of Law 2793: «Measures for the stabilisation of the public finance»). This opened a debate, which has not yet been concluded by a legislative act. Therefore policy options for a national insurance strategy are still open to investigation.

The case study has been developed for the Tuscany region (Amendola et al. 2000/a). Even if Tuscany is not among the most hazardous regions with respects to seismic activities, the case study is quite representative for the methodological approach, since all the results have been normalised to reference monetary units as described in the following. The choice of the region was determined by the fact that the Institute for Research on Seismic Risk of the Italian National Research Council made models and data from a previous study available (Petrini et al. 1995). These have been incorporated in a Monte Carlo generator of seismic events, which simulates occurrence of earthquakes affecting the region, calculates attenuation according to the geological characteristics, and finally determines the acceleration at the ground in each municipality. The IIASA spatial-dynamic, stochastic optimisation model has been customised to explicitly incorporate the vulnerability of the built environment, with data on number and types of buildings in each municipality of the region. In a first phase the study has been limited to the study of different policy options for an insurance program. In a next working phase the interplay between investments in physical mitigation (retrofitting) and risk-sharing measures should be investigated.

## 2. An overview of the model

The general model has been described in the quoted above references. In this case study, the Tuscany region has been subdivided into $M \approx 300$ sub-regions, which corresponds to the number of its municipalities. For each municipality $j = 1, 2, ..., M$, number and types of buildings (and therefore their vulnerability), and number of built cubic meters are available. These represent the so-called estimate $W_j$ of the property values or "wealth" of the municipality $j$. Simulated in time and space, earthquakes $\omega_0, ..., \omega_t$ may occur at different municipalities, inside or outside the region, have random magnitudes and, therefore, affect a random number of municipalities. From data and models in the Petrini et al. report, a catastrophe generator has been created, based on the Gütenberg – Richter law and on the attenuation characteristics of the region (for example, see Figure 1, Aniello et al. 2000, 2000/a). This enables the generator to calculate intensities and accelerations in each municipality. Of course, the generator could be easily adapted to incorporate different kinds of distributions, non-poissonian catastrophic processes, as well as micro-zoning within a municipality.

Fig. 1: Earthquake generator[2]


The municipalities affected at time t are indicated by a subset $\mathcal{E}_t$ of municipalities $j$, $j = 1,2,...,M$. Petrini's vulnerability relations between accelerations and losses according to the type (masonry or reinforced concrete), age and maintenance of the buildings are used to estimate the number of cubic meters of destroyed properties. The economic loss of destroyed cubic meters of a building is defined as the cost for their reconstruction. Then it is possible to be independent of contingent pricing by considering the cost of reconstruction per cubic meters to be the *monetary unit*. In this way the simulation of time histories for possible earthquakes in the region produces the sets of economical losses, and enables the design of an insurance programme.[3]

In its early version the Italian Design of Law 2793 (1998), to reduce the impact of natural disasters on the governmental budget, included in its Article 31bis provisions for an insurance programme against all natural hazards. It was intended not to make this insurance mandatory, but to make mandatory the extension of a fire insurance policy to all natural hazards, in a way similar to the French system quoted. In addition to tax incentives for such an insurance, it stipulated a maximum exclusion layer of 25%, the creation of a pool of insurance companies with a reserve fund corresponding to the annual average government payment for compensating losses (with some forms of state guarantee to be specified further), and linking of the premium to the premium for fire policy, rather than to the risks of a specific municipality. This article was withdrawn, and later proposals are still subject of discussion.

Starting from these principles, the case study intends to demonstrate how the model analyses and offers the decision-makers different policy options.

Let us assume that an insurance company (this might be a pool of companies or the government itself acting as an insurer) covers a fraction $q$ ($q = 0.75$ as in assumptions all owners buy the insurance) of earthquake losses. The rest, according to the Design of Law, would be

---

[2] Unpublished and work still in progress at IIASA by S. Baranov, B. Digas.

[3] It would also be possible to determine in which way preventive retrofitting could decrease the losses: this is easily done by a consequent decrease of the vulnerability indices in the loss model. In this way it would be possible to study the interplay between structural measures and risk-sharing for an integrated risk management approach, and to design an insurance system linked to incentives for retrofitting of the built environonment

compensated by the state. The state would also be exposed to feed the reserve funds in case of excessive losses. This would in any case allow the government to save money with respect to usually paid compensations and to use the save for prevention measures for public infrastructures and cultural heritage.

The company has an initial fund or a risk reserve $R^0$, which in general is characterised by a random variable dependent on past catastrophic events. To analyze necessary $R^0$, we can set different values, e.g., $R^0 = 0$. Assume that the time span consists of $t = 0,1...,T$, for example, $T = 50$ time intervals. The risk reserve $R^t$ of the company at time $t$ is calculated according to the following formula for $t = 0,1,...,T$:

$$R^t = R^{t-1} + \sum_{j=1}^{m} \pi_j^t - \sum_{j \in \mathcal{E}_t} L_j^t(\omega_t) q , \tag{1}$$

where $R^0$ is the initial risk reserve, $q$ is the coverage of the company in a municipality $j$ at time $t$, $\pi_j^t$ is the premium from the municipality $j$. Assume that $L_j^t(\omega_t)$ is the loss (damage) at $j$ caused by the simulated catastrophic event $\omega_t$ at time $t$. The value $L_j^t(\omega_t)$ depends on the event $\omega_t$ [4] and the type of properties in $j$. The analytical structure of the probability distribution of the random variable $R^t$ is intractable, therefore, the methodology relies on Monte Carlo simulation.

Usual actuarial approaches calculate their premiums with respect to the loss expectations. Therefore this study considered two policy options based on similar principles:

Premiums based on the average damage over all municipalities (solidarity principle, bringing less exposed locations to pay premiums equal to more severely exposed ones, as in the spirit of the proposed insurance programme)

Location-specific premiums based on average damage in the particular municipality (risk-based).

However, stochastic optimisation allows the analysis of different criteria and takes into account location specific, dependent risks. As an example, a third policy option has been considered:
Premiums calculated in a way that equalises in a fair manner the risk of instability for the insurance company and the risk of premium overpayment for exposed persons.
For Option 3 in this study, the insurer maximises his profits taking into account the risks of his insolvency under the constraint on 'fair' premiums. 'Fair' premiums are defined in the following sense. Let municipality $j$ face losses (damages) $L_j^t$. Individuals from this municipality receive a compensation $L_j^t q$ from the company when such a loss occurs. If $W_j^0$ is the initial wealth (property), then municipality-$j$'s wealth at time $t$ is

$$W_j^t = W_j^{t-1} + L_j^t q - \pi_j^t . \tag{2}$$

---

[4] In a general model it may also depend on time dependent mitigation measures or deterioration of the built environment.

It is assumed that individuals (municipalities) maximise their wealth according to the distribution of cases when $v_j^t < 0$, where $v_j^t = L_j^t q - \pi_j^t(\omega^t)$. Therefore, the 'fair' (optimal) vector of premiums $\pi^t = (\pi_1^t, ..., \pi_m^t)$ will guarantee the given level of stability for the insurer by minimising both the risk of his insolvency and the risk of overpayments for municipalities.

Thus, in a risk based or market approach, the choice of premiums reflects a certain balance between insurance demand and supply and creates additional incentives for insurance, otherwise higher premiums may decrease profits by decreasing the number of municipalities able to pay these premiums.

## 3. Numerical experiments

For *Option 1*, where the burden of losses is equally distributed over the population, the annual premium is equal to the flat rate of 0.02 monetary units (m. u.) per cubic meter of building (in per cent term, i.e., $\pi_j^t \times 100$, $j = 1, ..., m$, $t \in [0, T]$).

For *Option 2,* Fig. 2 shows the distribution of municipality-specific premiums based on average damage in each municipality (or according to the municipality-specific risk exposure factor). There is a prevailing number of municipalities (about 220) that have to pay 0.02-0.03 m. u., which is close to the flat rate of 0.02, as in Option 1. About 20 municipalities are at no risk at all (0 rate). Municipalities more exposed to the risk, have to pay 0.04 and higher rates (more than 50 municipalities).



Fig. 2. Distribution of municipality-specific premiums (per $m^3$ building volume at municipality, in per cent terms).

Fig. 3 shows the distribution of the insurers' reserve at premiums of Options 1 and 2. The reserve is cumulated over a 50 year time span. The volume of capital is defined by the horizontal axis. The probability of insolvency (when the risk reserve accumulated until the occurrence of the catastrophe is not enough to compensate incurred losses) is indicated on the

right-hand ordinate axis. As is seen, there is a rather high probability of 'small' insolvency (values –90, -40 occurred 190 and 90 times out of 500 simulations). High solvency (more than 500 m. u.) occurred in about 10 per cent of the simulations. The insolvency would represent the cost to the government in guaranteeing the reserve fund.

Fig. 4 shows the distribution of premiums for *Option 3*. According to this principle, most of the municipalities (190) have to pay close to the flat rate of 0.02-0.03 m. u. per cubic meter of building. Rates of 0.04 and higher have to be paid by about 100 municipalities. In this case the highest premium rate is 0.5, which, in comparison to the highest rate of 1.2 of Option 2, is much lower. The distribution of the insurer's reserve in Fig. 5 indicates also the improvement of the insurer's stability: the frequency of insolvency is reduced to 3 out of 500 performed simulations.



Fig. 3. Distribution of insurer's reserve for Options 1 and 2 (in thousands monetary units over 50 years).

Fig. 4. 'Fair' premiums according to model (1)-(2), or gainer-looser equilibrium (per $m^3$ building volume at municipality, in per cent terms).



Fig. 5. Distribution of insurer's reserve for Option 3 (in thousands monetary units over 50 years).

Fig. 6. Comparison of Options: flat, municipality-specific, and 'fair' premiums.

Fig. 6 is very illustrative. For each municipality it shows the optional premiums to be paid: the flat premium rate of 0.02, the Option 2 municipality-specific rate, and the 'fair' premium of Option 3. Many municipalities in all three options have to pay the premium rate, which is about the flat rate (0.015-0.03). For quite a number of municipalities in Options 2 and 3, the rate significantly exceeds the flat rate. Options 2 and 3, therefore, identify the municipalities, which are most exposed to the risks. For these municipalities special attention should be given as to whether they are able to pay such high risks. The model here incorporates average wealth (households' income) in municipalities, which can be regarded as additional constraint on the municipalities 'solvency' (overpayments). Option 3 allows to take such individual constrains into account and work out the premium rate optimal both for insurer and for municipalities.

## 4.   Conclusions

The case study based on a comprehensive geographically distributed data set has demonstrated the ability of the methodology developed at IIASA to analyse and compare different policy options for risk sharing in the case study region. The methodology is able to incorporate different kinds of hazard and vulnerability models, and to deal with various kinds of dependencies.

Future work should

include probability distributions for vulnerability, instead of point value, as in the present study;

investigate the trade-off between structural mitigation measures and insurance strategies for an integrated risk management;

include a people's behaviour model with respects their willingness to make use of possible incentives to reduce vulnerability, and/or to buy insurance. In this case the live- savings aspects of retrofitting should also be considered;

introduce dynamics of reconstruction, and superimposing seismic crises.

32

# References

1. Amendola, A., Ermoliev, Y., Ermolieva, T., Gitis, V., Koff, G., Bayer-Linnerooth, J., 2000. A Systems Approach to Modeling Catastrophic Risk and Insurability. Natural Hazards, 21: 2/3, pp.381-393

2. Amendola, A, Ermoliev, Y., and Ermolieva, T.Y., 2000/a. Earthquake Risk Management: A Case Study for an Italian Region. Web proceedings of the EuroConference on Global Change and Catastrophe Risk Management: Earthquakes Risks in Europe. http://www.iiasa.ac.at/Research/RMP/july2000/

3. Cummins, J.D., and Doherty, N., 1996, Can insurers pay for the "Big One"? Measuring capacity of an Insurance market to respond to catastrophic losses. Wharton Risk Management and Decision Processes Center, University of Pennsylvania, W.P.

4. Digas, B., 1998, Generators of Seismic Events and Losses: Scenario-based Insurance Optimization. IIASA, Interim Report, IR-98-079.

5. Ermolieva, T., Ermoliev, Y., and Norkin, V., 1997, On the role of advanced modeling in managing catastrophic risks. in Drottz-Sjöberg, B.-M. (ed.), Proc. New Risk Frontiers: Conference for the 10th Anniversary of the Society for Risk Analysis - Europe, The Center for Risk Research, Stockholm, pp.68–74.

6. Ermolieva, T., Ermoliev,Y., and Norkin, V., 1997, Spatial Stochastic Model for Optimization Capacity of Insurance Networks under Dependent Catastrophic Risks: Numerical Experiments. IIASA, Interim Report, IR-97-028.

7. Ermoliev, Y., Ermolieva, T., MacDonald, G., and Amendola, A., 2000, A systems approach to catastrophe management. *European Journal of Operational Research* ,122:452-460.

8. Gilber, C. and Gouy, C., 1998, Flood Management in France, in Rosenthal U. and Hart, P't eds. Flood Response and Crisis Management in Western Europe: A Comparative Analysis. Springer. Berlin.

9. Kunreuther, H., and Roth, R., 1998, Paying the Price: The Status and Role of Insurance Against Natural Disasters in the United States, Joseph Henry Press, Washington, D.C

10. Linnerooth- Bayer, J. and Amendola, A., 2000, Proceedings of workshop on Mitigation of Seismic Risk-Support to Recently Affected European Countries, November 27-28, Belgirate (VB), Italy.

11. Linnerooth- Bayer, J., and Amendola, A., 2000/a. Global Change, Catastrophic Risk and Loss Spreading. The GENEVA PAPERS on Risk and Insurance. 25:2, 203-219.

12. Linnerooth-Bayer, J., and Quijano, S., 2000, Loss Sharing: A Study of Recent Major Earthquakes and Floods. Conference Global Change and Earthquakes Risks. http://www.iiasa.ac.at/Research/RMP/july2000/

13. Pasterick, E. T., 1998, The National Flood Insurance Program, Paying the Price: The Status and Role of Insurance Against Natural Disasters in the United States, In Kunreuther, H. & Roth, R. eds., Joseph Henry Press: Washington, D.C., pp: 125-155

14. Petrini V. et al., 1995, Pericolosita' sismica e prime valutazioni di rischio in Toscana. CNR/IRRS, Milan.

# Part II: Risk Management and Modeling Techniques

# The case for innovative nuclear reactor and fuel cycle systems

*Viktor Mourogov, Vladimir Kagramanian*

## Abstract

This paper is divided into three sections. The first one describes the changes taking place in the market for NPPs and the resulting need for innovative reactors and fuel cycles. The second one outlines the range of innovative approaches that have been identified as well as requirements for success. Finally, the paper addresses the need for international R&D cooperation, international initiatives that are already underway, and the role of the IAEA.

# 1. Introduction

The potential market for nuclear power plants (NPPs) is changing. To capture a significant part of the future potential market, NPP technology will therefore have to change too. If new reactors and fuel cycles are not developed so as to be economically competitive in the changing market and attract public acceptance in terms of safety, spent fuel and waste management, and proliferation resistance, then other power technologies will fill the gap.

# 2. The need for innovative reactors and fuel cycles

Nuclear power is now at a crossroads, with no consensus concerning its future. A number of factors currently favor nuclear power, including:

- Nuclear energy has grown in only 50 years from a new scientific development to a major part of the energy mix in several of the 32 countries now using nuclear power. The majority of currently operating NPPS perform well, both in terms of economics and safety (IAEA, 2000a; IAEA, 2000b).

- In addition to producing electricity, nuclear power helps meet other important national goals in many countries, including energy independence, clean air, and reduced emissions of greenhouse gases (GHGs).

- Nuclear power is currently the only mature non-carbon electricity generation technology that can significantly contribute to the long-term global sustainable energy mix. Recent scenarios developed by the Intergovernmental Panel on Climate Change (IPCC) foresee a significant potential for the nuclear energy growth - from its current 6% of primary energy to between 10% and 30% by 2100 (IPCC, 2000). This implies an increase in global nuclear power capacity from current levels of 350 GW(e) to 2,000–5,000 GW(e) by 2050 (Kagramanian et al, 2000).

Against this background, one could reasonably expect nuclear power generation to be in the midst of a rising trend extending out into the foreseeable future. Yet, that is not the case. Figure 1 shows average annual new capacity additions, worldwide, for 5-year periods from 1971 to today. As is evident in the figure, new nuclear additions have been decreasing since the middle of 1980s and are now an almost vanishingly small part of total capacity additions. For example, in 2000 six new power reactors (India-3, Pakistan-1, Brazil-1, and the Czech Republic-1) with a total capacity of 3GW(e) were connected to the grid. This equals only 3% of estimated total global annual capacity additions in 2000. With nuclear power's estimated share of global electricity production holding steady at about 16% in 2000, nuclear's share of new capacity was less than one fifth its share of electricity production. Quite simply, over the last 15 years, nuclear power has been losing market share badly in a growing world electricity capacity market.

**World Electricity Capacity Additions History**

Fig.1.

Looking at the developed regions that have principally supported nuclear technology, we find a relatively greater contribution of nuclear power, but not necessarily brighter prospects for the future. In North America, Fig.2 shows that capacity additions declined in the first part of the period presented, driven by energy conservation, economic slowdown, and some overcapacity. Nuclear power maintained a significant share, however, through the 1980s. Now capacity additions have started to grow again, but net additions for nuclear power have dropped first to zero (1991-95) and then below zero (1996-97).



**NA Annual Electricity Capacity Additions**

Fig 2.

Fig. 3 shows a similar situation in Western Europe. New additions have begun to rise, but it is clear that nuclear power, despite its substantial role in Western Europe, will not be part of that rise for at least the next decade and quite probably more.



Fig.3.

Fig. 4 presents the history for Eastern Europe. Here the situation is dominated by economic recession. But whatever the cause, the picture for nuclear is not promising.



Fig.4.

Fig. 5 describes the situation in most of the world's developing regions other than China. The figure includes Latin America, Africa, the Middle East, South and Southeast Asia, and the Pacific.



Fig.5

Fig. 5 gives no indication of nuclear power expanding significantly beyond the developed countries where it was born, to the developing regions which will eventually dominate global energy growth.

Fig. 6 presents developments in the Far East, which we usually emphasize as the best current and future market for nuclear. Here, the picture is brightest. Overall capacity additions are rising fastest and nuclear's share is the most distinctive. But even here, nuclear capacity additions are basically staying about constant, not growing with the market. Even here, the result is a diminishing market share in a growing market.



Fig.6.

Overall, it's a quite sobering picture. The rate of new electricity generating capacity additions has started to grow in the last 10 years after a prior slowdown. But nuclear additions are hardly keeping pace. Thus we have a growing market, but a small and shrinking nuclear share.

For the immediate future - the next one or two decades - things do not look much better. In Fig. 7, which shows the IAEA's projected capacity additions through 2020, the overall picture for the global electricity sector, in both the high and low scenarios, is quite healthy. But not so for nuclear. In the low scenario net global nuclear capacity additions between 2011 and 2020 are negative. In the high scenario, they are so small as to be unnoticeable for the next decade.



**World Electricity Capacity Additions: IAEA Projections**

Fig.7

Other projections, for example from the OECD/NEA and US DOE, present a similar picture. Nuclear power has been losing market share in a growing market, and according do these projections it will continue to do so - the only question is how fast.

Part of the reason for the bearish outlook in Fig. 7 is that in several western countries, including some with the most expertise in the nuclear field, there is strong current political opposition to nuclear power. Opponents emphasize concerns about nuclear waste, safety and non-proliferation. The present governments of Germany, the Netherlands, Belgium and Sweden - countries that have relied heavily on nuclear power -intend to phase it out. Political opposition now also extends, in some cases, to nuclear power's inclusion among the sustainable energy technologies eligible for Clean Development Mechanism (COM) projects under the UN Framework Convention on Climate Change (UNFCCC).

In North America and Western Europe, nuclear power capacity additions also faces economic challenges arising from electricity market deregulation, plus improvements in particularly gas-fired power generation. The situation has been aggravated by low electricity demand growth over the past decade.

In developing countries, the main issues limiting nuclear expansion are a lack of expertise in nuclear technology and its safety culture, a lack of adequate infrastructure, and of course economic and financing issues.

To address these challenges, the nuclear community has made progress on a number of fronts.

Fig. 8 shows the record of improvement in the aggregate availability of nuclear power plants over the last decade. The trend is significant and substantial. It is the equivalent of 28 GW(e) of new capacity - highly cost-effective new capacity. However, this route to success has its limits. Realistically the trend in Figure 8 will have to level out somewhere around 85% or 90%.

**World Energy Availability Factor by Year**

Fig.8

Lifetime extension of existing NPPs is another important accomplishment. But again this route to success has its limits. Figure 9 shows the age distribution of today's operating reactors. In the middle is the "baby boom" from the 1970s and 80s. For these reactors, lifetime extensions promise substantial capacity extensions. But after the "baby boom" reactors, there is not much in the pipeline. The reality is that the impact of lifetime extensions can only diminish.



**Number of Reactors in Operation by Age**

Fig.9

New institutional measures have been organized in the areas of safety, waste and non-proliferation, which have helped to significantly reduce the risks of severe accidents and of proliferation of fissile materials within existing nuclear power system. The impressive health and environmental advantages of nuclear power, relative to alternative electricity generation options, have been analyzed and advertised. However, to the extent that the corresponding

41

disadvantages of fossil-fuel alternatives are not internalized in capital and operating costs, nuclear's health and environmental benefits have little impact on investors in new generating capacity. They prefer the cheapest option, and most often it is not nuclear.

The industry continues work on evolutionary new designs to improve NPP performance even further. New designs for advanced LWR reactors may well have a market in a limited number of Asian countries that do not have access to cheap gas from pipelines, such as the Republic of Korea and Japan. Elsewhere it will be extremely difficult, if not impossible, for new evolutionary advanced nuclear power plants to compete economically with gas-fired plants in the absence of new governmental policies.



Fig. 10



Fig. 11

1. Countries with scheduled plutonium utilization programs (Belgium, France, Germany, Japan, Switzerland, and others)
2. Countries with no scheduled plutonium utilization programs (Russia, UK, US, and others)

In short, the nuclear community has accomplished much, but evidently not enough to turn around the trends shown in Fig. 1-7.

We believe that the main reason for this stagnation is that the nuclear community continues to rely on nuclear technology developed in the 1950s. This technology has its roots in military applications and does not easily lend itself to the features that might be characteristic of new designs qualifying as "inherently safe." We believe this technology has now reached its limits from an economic point of view. To the extent we continue to rely on this technology, any substantial expansion of nuclear power will therefore necessarily depend increasingly on improved human performance, active safety systems, and institutional and organizational measures to reduce accident and proliferation risks. Such a route essentially increases "overhead" without improving the underlying technology. It is unlikely to increase economic competitiveness, thereby leaving any expansion of nuclear power to be driven largely by external factors, such as fossil fuel prices and environmental taxes.

If the nuclear power sector is to increase its role, it cannot simply continue to do what it has been doing and expect that factors outside its control, such as fossil fuel prices or environmental taxes, will change to make nuclear power's prospects more favorable. To reach a different outcome than that indicated by current trends, something must be done within the nuclear community to generate new technological solutions. The challenge is to look to the future, to identify what innovations and new directions - that build upon and make good use of existing expertise and accomplishments - are most promising for helping nuclear power capture a growing share of a growing market. More of the same won't do. The industry must look to the future and must be innovative.

## 3. Objectives and Approaches to Innovation

Innovation is the foundation of technological and social development, and an essential feature of all commercially viable technologies. Within the nuclear power industry, innovation can take place on several levels:

- Improvements in maintenance, operations and other practices for existing commercial facilities.
- Incremental improvements in existing commercial reactor designs.
- New and advanced changes in design and operation involving major departures from current commercial designs.

Continuing improvements in the first two categories are critical, particularly for the next two decades, to keep nuclear power alive while innovations falling in the third category are developed. It is these innovations in the third category – involving major departures from current commercial designs – that are the focus of the balance of this paper and are the innovations referred in the papers title. These will be essential for a new generation of reactors and fuel cycle(s) that can form the basis for a large-scale worldwide expansion of nuclear power.

The need for innovative R&D has been recognised by the nuclear industry and by those countries that believe in the overall benefits, viability and importance of nuclear power for the long term. Currently, R&D on innovative nuclear fuel cycle and reactor concepts is being performed in a number of countries, including Argentina, Canada, China, France, India, Italy, Japan, the Republic of Korea, Russia, South Africa, and the USA. The USA in particular embarked on a Nuclear Energy Research Initiative in 1999 to develop advanced reactor and fuel

cycle concepts and promote scientific breakthroughs to help overcome the current obstacles to nuclear expansion (Majumdar et al, 2000; Mourogov and Kupitz, 1998; Magwood, 1999).

The IAEA has been active in assisting Member States in developing advanced reactor and fuel cycle technologies within the framework of established International Working Groups. We have found a wide diversity of requirements, or targets, used to guide the development of innovative reactor and fuel cycle concepts. Not surprisingly, the result is a wide diversity of reactor and fuel cycle concepts. The following are several examples of advanced reactor and fuel cycle concepts, and of the targets guiding their design.

- Small modular reactors (HTR) with once-through fuel cycles. Targets: competitiveness; minimum financial risk; 2-3 year construction period; increased safety (no release of radioactivity from fuel); potential for both electric and non-electric applications, and for use in areas with or without grids.
- Fast reactor (FR) and closed fuel cycle concepts. Targets: competitiveness; deterministic safety; effective use of uranium; transmutation of long-lived radioactive and toxic actinides and fission products; proliferation resistance.
- Accelerator driven systems (ADS). Targets: transmutation of long-lived radioactive and toxic actinides and fission products.
- Thorium based reactor and fuel cycle systems. Targets: competitiveness; expanded resource base; proliferation resistance.

In addition to differences in the types of requirements and targets that guide these efforts, there are also differences within each area of concern, e.g., economics, safety, waste management, and non-proliferation. In connection with economic competitiveness, for example, opinions differ as to whether or not assessments of competitiveness should take into account the potential introduction of CC>2 taxes and increases in fossil fuel prices. In such crucial areas as safety, waste management, non-proliferation, resource consumption, and the application for which the new reactor is principally intended, an even wider diversity of requirements and targets exists. For safety, for example, some believe that today's advanced LWRs are sufficiently safe for large-scale deployment because they are neighbor-friendly (i.e., the probability of a significant release of radioactivity off-site even in the case of severe accident is negligible). Others insist that the public will accept large-scale nuclear energy deployment only of new reactors, like modular HTRs or innovative FRs, with deterministic safety features (i.e., no significant release of radioactivity at all).

Concerning waste management, some believe that direct underground disposal of small volumes of spent fuel is a sufficiently safe back-end option. All that is needed to assure public acceptance is a practical demonstration. Others insist that only the elimination of long lived hazardous nuclides by burning or transmuting them in fast reactors or ADS would attract public support for large-scale nuclear energy development. However, even among advocates of the latter view, there remain differences of opinion as to which hazardous elements should be reduced, and to what extent.

Concerning non-proliferation, some believe that the continuing development of the safeguard system along its present trajectory will be sufficient to safeguard effectively any new reactor and fuel cycle. Others propose special "proliferation resistant" reactor and fuel cycle concepts (using new types of fuel, new reprocessing technologies without plutonium extraction, new FR concepts, etc.) that rely increasingly on technical measures against possible proliferation. There remains no consensus among researchers about how to measure proliferation resistance, or about the extent to which reliance on technical measures should be increased.

## 4. International Cooperation and the Role of the IAEA

Both the importance of national activities on innovative reactor and fuel cycle systems, and the desirability of coordinating them internationally, have been acknowledged at several meetings held under the auspices of the IAEA. The IAEA Scientific Forum (September 1999), the Advisory Group Meeting on Development of a Strategic Plan for an International Research and Development Project on Nuclear Fuel Cycles and Power Plants (October 1999), the Industry Forum (January, 2000). These meetings have recommended that the Agency take steps to help assess the potential for, and facilitate the exchange of information on, innovative nuclear reactors and fuel cycles.

The Agency has for more then forty years served "to accelerate and enlarge the contribution of atomic energy to peace, health and prosperity throughout the world and to ensure, so far as it is able, that assistance provided by it or at its request or under its supervision or control is not used in such a way as to further any military purpose" (IAEA, 1956). The Agency's on-going activities on innovative reactor and fuel cycle technologies include the activities of several international working groups and co-ordinated research projects. Furthering international co-operation has always been understood to constitute one of the main areas of activities where the Agency can be of benefit to its Member States.

Member State interest in new future applications of nuclear power has been increasing in recent years, culminating in a resolution by the 2000 General Conference supporting an IAEA initiative on innovative nuclear power reactor technology development. In response to this resolution, the IAEA has recently launched an "International Project on Innovative Nuclear Reactors and Fuel Cycles", INPRO. At a meeting of senior officials from Member States and international organizations in Vienna in November 2000, the objectives and conditions of this project were discussed and the Terms of Reference for INPRO were finalized. The objectives of INPRO, as defined in the Terms of Reference, are:

- to help to ensure that nuclear energy is available to contribute in fulfilling, in a sustainable manner, energy needs in the 21st century;
- to bring together all interested Member States, both technology holders and technology users, to consider jointly the international and national actions required to achieve desired innovations in nuclear reactors and fuel cycles that use sound and economically competitive technology, are based - to the extent possible - on systems with inherent safety features and minimise the risk of proliferation and the impact on the environment;
- to create a process that involves all relevant stakeholders that will have an impact on, draw from, and complement the activities of existing institutions, as well as ongoing initiatives at the national and international level.

The Project will be implemented in two phases. Phase I will be initiated in the end of the April 2001. In the first phase, work will proceed in five subject areas recognised as important for the future development of nuclear energy technology. The five subject areas are: Resources, Demand and Economics; Safety; Spent Fuel and Waste; Non-proliferation; Environment. The work will include selection of criteria and development of methodologies and guidelines for the comparison of different concepts and approaches and determination of user requirements in the subject areas.

Upon successful completion of the first phase, taking into account advice from the Steering Committee, and with the approval of participating Member States, a second phase of INPRO may be initiated. Drawing on the results from the first phase, it will be directed to:

- examining in the context of available technologies the feasibility of commencing an international project;
- identifying technologies, which might be appropriate for implementation by Member States of such an international project.

The INPRO will take advantage of the IAEA's global membership in addressing the global dimensions of nuclear energy issues (e.g., technology, safety and safeguards) and the fact that future nuclear energy demand growth will occur mostly in developing countries. The INPRO will complement existing international efforts such as the three-Agency studies on Innovative Reactor Development by OECD/IEA, OECD/NEA and IAEA and the US initiated Generation IV International Forum. These have objectives related to those of the INPRO, but with some differences in either scope or participation. Good co-ordination and co-operation with other initiatives have to be pursued in a complementary manner to take full advantage of potential mutual benefits.

What we can contribute, in addition to our global membership, which includes both new technology-user countries, old technology-supplier countries, and everything in between, is also our unique collective expertise in safeguards, safety and technology. A really successful innovative design or fuel cycle is unlikely to be one that is exclusively focused on just the economics, or just safety, or just proliferation resistance. It will have to represent a distinctive improvement on many fronts all at once, and that requires the sort of exchanges among specialties and cross-fertilization of ideas where the IAEA can be of assistance. Moreover, for a new design to become successful on a large scale, and make a substantial contribution to future energy supplies, it must succeed at a global level. A global perspective and understanding is, again, something that the IAEA is uniquely positioned to offer, given its membership, mission, and experience.

## References

1.  IAEA, 1956, Statute of the International Atomic Energy Agency, Article II.

2.  IAEA, 2000a, Nuclear Technology Review 2000, IAEA GOV/INF/2000/5/Part 1 IAEA (2000b), IAEA PRIS Database as of 01 January 2000.

3.  IPCC, 2000, Working Group III: Mitigation of Climate Change, Special Report on Emission Scenarios, Intergovernmental Panel on Climate Change.

4.  Kagramanian, V.S., Kononov, S.L., and Rogner, H.H., 2000, Nuclear Energy in the New IPCC Emission Scenarios, IAEA Bulletin, Vol. 42, No. 2.

5.  Majumdar, D., Kupitz, J., Rogner, H.-H., Shea, T., Niehaus, F., and Fukuda, K., 2000, The need for innovation, IAEA Bulletin, Vol. 42, No. 2.

6.  Mourogov, V.M., and Kupitz, J., 1998, Nuclear Energy Issues and the Role of Small and Medium Sized Reactors, 23rd Annual Symposium of Uranium Institute, London, U.K., 9-11 Sept. 1998.

7.  Magwood, W. D., 1999, US Department of Energy Generation IV Nuclear Power Systems, ANS Meeting, November 16, 1999.

8.  Fukuda, K. , Bonne, A., and Mourogov, V.M., 1999, Global View on Nuclear Fuel Cycle -Challenges for the 21st Century, Proceedings of the International Conference on Future Nuclear Systems Global'99 "Nuclear Technology -Bridging the Millennia" August 29 -September 3, 1999 Jackson Hole, Wyoming, USA.

# Soft computering approaches in fault diagnosis system and risk management

*Józef Korbicz*

## Abstract

The paper deals with the Group Method of Data Handling that belongs to a class of evolutionary algorithms. The theory of the network of the GMDH type, which enables to solve some problems that cause the restrictions of the classic solutions of the artificial neural networks, is outlined. Some perspectives of the development of the theory of the GMDH type neuron networks are shown.

# 1. Introduction

The very well known solutions of the neuron networks can be applied in different fields of science and technology due to such features as simplicity of implementation, good approximation of the built-up systems and the possibility of convenient formation of equipment applications. An increase in the complexity of the examined objects and increasing requirements for efficiency and reliability of the used analytical tools stimulate the search for new solutions. The research connected with integration of the artificial neuron networks and the other well known methods of artificial intelligence are particularly interesting and promising.

In particular the tasks of modeling the dependence between the output and input quantities in complex objects show weak points of the classic works on the neuron networks. The applied inductive methods that use only empirical data for the synthesis of the general rules of the model construction turn out to be less effective when they are used in learning the arbitrary defined network structure. In case of a complicated model and data with noises this factor has an impact on the accuracy of the achieved results. The concept of extending the learning process to the typology of connections among neurons should be taken into consideration. The Group Method of Data Handling GMDH that belongs to a class of evolutionary algorithms can be useful for this purpose. The definition of the proposed extension method is flexible and besides the structural and parametrical network optimalization it also allows to extend the scope of applications to the dynamic systems.

# 2. GMDH algorithm

The neuron networks (Korbicz et al. 1994, Tadeusiewicz 1993) are the most efficient analytical tools used in case when the parameters and the model structure are not known. Due to their applications the unknown structure of the mathematical model need not to be defined. There is, however, another problem that results from the necessity to define the network structure. LMS algorithm allows for a given task the parametrical optimalization of the unlimited number of forms of the modeling functions. Analogically the artificial neuron network can be trained for many various typologies of the internal layers of neurons. Therefore, despite the well known effectiveness of the presented class of mathematical apparatuses they introduce to the solution the method error difficult to evaluate.

In case of dependency test between the input and output signals of the complex objects, this error can have an essential impact on the practical use of the achieved results. The search for new solutions in modeling algorithms based on processing empirical data is justifiable (Kohonen 1984, Hecht-Nielsen 1991). Group method of data handling – GMDH is the example of such a solution.

## 2.1 General assumptions

The group method of data handling was developed in the late sixties by A.G. Ivakhnenko from the Ukrainian Academy of Science. It became well known after the foreign publications in seventies (Ivakhnenko 1971). Ivakhnenko developed the GMDH algorithm to predict accurately the evolution of fish population in rivers and oceans. The main idea of this method was the synthesis of the polynomial model. Because of the innovative integration of a few concepts of the structural and parametrical optimalization, Ivakhnenko's polynomial, the result of the GMDH procedure, turned out to be a model ensuring unusual accuracy and practical use.

The basic assumption of the presented algorithm was the resignation from the deductive approach based on the knowledge of engineers and experts. The other essential element was the idea of polynomial evolution from the elementary to optimal structure through the selection of

various combinations of the simple partial models. In the majority of works they are polynomials of second degree of two variables. In this concept the degree of resulting polynomial doubles in any algorithm stage, taking into account that in every subsequent iteration polynomial functions construed in previous iteration are arguments of the elementary models. The optimal values of the constant parameters are calculated by means of the least square method. In the described way the resulting structure of high degree of compilation and optimalized parameters is possible to achieve in a few steps. The application of Kurta Gödl's theory (Nagel i Newman, 1966) allowed Ivakhnenko to propose the criterion of detecting the optimal complexity of such a developed model.

The assumption introduced in the group method of data handling which defined the partial polynomial as approximation **N-*th*** degree (most often N=2) of Maclaurin's series defined for the polynomial function, known as Kolmogorov-Gabor polynomial (Ivakhnenko 1971) was very important.

$$y = a_0 + \sum_{i=1}^{N} a_i \cdot u_i + \sum_{i=1}^{N}\sum_{j=1}^{N} a_{ij} \cdot u_i \cdot u_j + \sum_{i=1}^{N}\sum_{j=1}^{N}\sum_{k=1}^{N} a_{ijk} \cdot u_i \cdot u_j \cdot u_k + ... \qquad (1)$$

Two goals have been achieved this way. First – linear dependence of the elementary function from constant parameters was achieved, which allows for easy optimization e.g. by means of LMS algorithm. Secondly – function (1) can be written in a general form (for $N \geq 2$):

$$y = f_1(u_1) + f_2(u_2) + f_{12}(u_1, u_2) \qquad (2)$$

Therefore any argument of a partial function is formed by modification (functions $f_1$ and $f_2$) of two arguments from the earlier iteration of algorithm and their combination (function $f_{12}$). Dependencies $y=f_1(u_1)$ and $y=f_2(u_2)$ constitute a general record of the *mutation* data function, and $y=f_{12}(u_1,u_2)$ is a general record of *crossing* (exchange of information). Therefore, in the GMDH method you can find elements of the theory of *genetic algorithms*. Although Ivakhnenko did not directly consider these analogies, he underlined the importance of evolutionary mechanisms in his algorithm many times.

Many applications that appeared after publication of the assumptions of the GMDH algorithm confirmed its effectiveness and wide scope of applications. In early seventies Adaptronics Inc., an American firm used GMDH, among others, in quality control apparatus of materials used in aviation industry. The attempts of equipment implementation of the group method of data handling procedures showed the analogies to the other similar solutions in the fields of artificial neuron networks. (See Fig. 1).



Fig.1. An example of a scheme of hardware realization of the GMDH elementary polynomial function

The hardware realization of the partial polynomial function of two variables of a second degree presented in Fig.1 can be associated with a neuron model. The implementation of Ivakhnenko polynomial, formed from the connection of such component elements, can be considered as an artificial neuron network. This is why the Adaptronics Inc. classified its applications of the group method of data handling as neuron networks. Soon Ivakhnenko himself also started to use this interpretation of his work.

In the described way a new class of neuron networks was formed on the basis of the group method of data handling. New effective learning apparatus based on the principle of feed-forward was introduced while the general typology typical of the other solutions of this type was maintained. Thanks to this, the mechanisms of structural and parametrical optimization were introduced to the network of the GMDH type. This network grows and develops in the learning process until the development leads to the improvement of the performance effectiveness. Thanks to the integration of elements of a few mathematical apparatuses (such as e.g. neuron networks, genetic algorithms, the least squares' method, Gödl's theory, clusterization) the GMDH network theory is a flexible and open construction for various types of modification. The method of incorporating the mentioned additional algorithms in the process of network synthesis has an influence on the characteristic features of the described solution.

## 2.2 Neuron model

The neuron network is constructed through the connection of a given number of the elementary cells (neurons) processing signals. The scheme of a single neuron of the GMDH type presents Fig. 2.



Fig.2. The GMDH neuron structure

It is assumed that always at least two input signals $u_1$, $u_2$, ..., $u_k$ constitute the stimulation which results in the formation of the output signal $y$. Dependencies between the two group of signals describe the transition function $f$ defined by the general definition (3).

$$y = f(\mathbf{u}) = f(u_1, u_2, ..., u_k) \tag{3}$$

The GMDH algorithm allows much freedom in defining function $f$. In most cases it has a non-linear character. The precise definition of the transition function is not of great importance because the concept of the evolution of the network model from a simple to complex form is considered here. The selection and neurons optimalisation mechanism will lead to the optimal global solution regarding from what the initial form the evolution started. It is essential for $f$ to meet the conditions defined by dependence (2) only. Therefore, the processing and exchange of information between the elements of network will run more effectively, which will result in minimizing the risk of finding a solution that is locally optimal.

From the practical reasons the equation (3) should not be too complex because it would complicated the learning process and extend the time of its durability. The minimal processing

error for the resulting network (see chapter 2.3) will be defined more precisely when its complexity will be gradually increased in the subsequent layers. In a classical approach the transition function has a form of a polynomial. Most often it is the approximation of *N*-th degree of Kolmogorov-Gabor polynomial (1). The demands for simple form and compatibility with a form (2) are completed for *N*=2. This is why in the majority of works related to the GMDH network the following definition is applicable (4).

$$y = a_0 + a_1 \cdot u_1 + a_2 \cdot u_2 + a_{11} \cdot u_1^2 + a_{22} \cdot u_2^2 + a_{12} \cdot u_1 \cdot u_2 \qquad (4)$$

From a form of the definition (3) and (4) results the feature differentiating the GMDH type of network from the other solutions. Values $a_i$ play a role of the constant parameters rather than scales of connections, which can be found in the majority of the well-known neuron definitions. As a consequence, a specific algorithm of learning and the specified features of network are described further in this work.

A characteristic feature of the group method of data handling is training of every neuron separately before it is incorporated to the next layer (see chapter 2.3). In this way every GMDH neuron „tries" to become the network output modeling the searched dependence $y=f(\mathbf{u})$. This is the only criterion showing the process of network evolution in any stage of the realization of algorithm.

The approval of the definition (1) or its particular case (4) allows to simplify the learning process by means of the least square method. The subset of learning data is used for this purpose. If it includes *n* measuring vectors $(u_1, u_2, ..., u_m, y)$, then for each pair $u_i$, $u_j$ one can write the system of equations (Farlow 1984):

$$y = a_0 + a_1 \cdot u_{1i} + a_2 \cdot u_{1j} + a_{11} \cdot u_{1i}^2 + a_{22} \cdot u_{1j}^2 + a_{12} \cdot u_{1i} \cdot u_{1j} \qquad (5)$$

$$y = a_0 + a_1 \cdot u_{ni} + a_2 \cdot u_{nj} + a_{11} \cdot u_{ni}^2 + a_{22} \cdot u_{nj}^2 + a_{12} \cdot u_{ni} \cdot u_{nj}$$

that after recording in a matrix form has a form

$$\mathbf{Y} = \mathbf{U} \cdot \mathbf{A} \qquad (6)$$

The set of Gauss's equations is formed after bilateral multiplication (6) by the transposed vector of the input signals:

$$\mathbf{U}^T \cdot \mathbf{Y} = (\mathbf{U}^T \cdot \mathbf{U}) \cdot \mathbf{A} \qquad (7)$$

which solution (8) allows to find parameters **A** minimizing the mean-square error of equations (6) for the set of training samples.

$$\mathbf{A} = (\mathbf{U}^T \cdot \mathbf{U})^{-1} \cdot \mathbf{U}^T \cdot \mathbf{Y} \qquad (8)$$

The matrix **A** includes the optimized values of the coefficients of the transition function of the form (4), for which each neuron aproximates value of the output signals with a minimum mean square error. This procedure is repeated for each element of the network formed for all the combinations of the input signals $u_i$. The neuron parameters determined in training procedure are not subject to any changes in the process of further synthesis of network.

The advantage of the presented solution is a simple computation algorithm that gives good results even for small sets of measuring data. The drawback, however, is the necessity to perform many matrix calculations. This is why for many large sets of measuring data it is necessary to sort them out first by means of clusterization technique (Duran and Odell, 1974) or to use the alternative training rules.

The methods known from the classical solutions of the neuron networks can also be used instead of the proposed by Iwachnienko equations of regression to train the GMDH neurons. Such methods are for example Widrow-Hoff's rule (Pham and the others 1995), that dissolve the task of determining the parameters of transition function by means of tested recurrence rule. Some vectors of parameters and signals have a form respectively (9) and (10).

$$\mathbf{A} = \begin{bmatrix} a_0 & a_1 & a_2 & a_{11} & a_{22} & a_{12} \end{bmatrix}^T \qquad (9)$$

$$\mathbf{U} = \begin{bmatrix} 1 & u_i & u_j & u_i^2 & u_j^2 & u_i \cdot u_j \end{bmatrix}^T \tag{10}$$

The Widrow-Hoff's delta rule adapted to determine the coefficients **A** depends on recurrent use of dependence (11) for subset of training data.

$$\mathbf{A}_{k+1} = \mathbf{A}_k + \alpha \cdot \frac{\mathbf{U}_k}{|\mathbf{U}_k|^2} \cdot \left( y_k - \mathbf{A}_k^T \cdot \mathbf{U}_k \right) \tag{11}$$

where $y_k$ is a value of output signals in $k$-th measuring sample, and $\square$ is a training coefficient. In practice $\square$ is applied from the range of values (0.1, 1). The application of equation (11) causes gradual modification **A** leading to reduction of difference between the measured and approximated value of output signals. Every neuron is trained separately by the described method until when the mean-square error achieves minimum.

The majority of the other training algorithms known from the theory of neuron networks, such as rules of Hebb's, Kohonen's, Grossberg's and others (Korbicz i inni 1994, Tadeusiewicz 1993), can be applied in case of the GMDH network. It is important to remember about restrictions connected with them that require the necessity to determine the initial values of parameters, the risk of over training of the network and a long training process, etc.

The other problem that sometimes makes to withdraw from the classical solution of the group method of data handling is a form of the transition function. In majority of applications the polynomial model is sufficiently accurate due to the mechanisms of learning and neuron selections introduced by the GMDH and the evolution of the network structure. In some cases however (e.g. research of oscillation signals) the use of dependence (3) of the trigonometric functions (Iwachnienko and the others 1987) can be more effective. This is why although the GMDH algorithm belongs to the class of the inductive methods it is possible to use expert's knowledge to improve the accuracy of data processing in a construed network. The theory of GMDH network allows certain modifications at the stage of defining a single neuron if they are within the allowed scope i.e. they do not cause the modification of the other steps of the algorithm of the network synthesis.

## 2.3   Network synthesis

The concept of the synthesis of the artificial neuron network of the GMDH type is based on iterative processing of a defined sequence of operations leading to the evolution of the resulting structure. The process is completed once the optimal degree of complexity is achieved.

In the first iteration, the input layer of neurons described by transition functions (3) for all the combination of signals $u_i$ is formed. In a general case of network of $m$ inputs built up from $k$-input neurons ($m > k$) $s$ new elements are formed this way, where $s$ has a form (12).

$$s = \binom{m}{k} = \frac{m!}{k!(m-k)!} \tag{12}$$

The use of the multi-input neurons is not advisable because of the computation reasons because the value $s$ grows considerably, not causing any increase of the accuracy of the network processing. For the transition function defined by a form (4) the number of neurons construed in the input layers equals to the amount of the possible pairs $u_i$ :

$$s = \frac{m \cdot (m-1)}{2} \tag{13}$$

Then the first layer of the GMDH network has a structure presented in Fig. 2.3, and its component cells are described in the transition functions respectively (14).

Fig. 3. Input layer of the GMDH neuron network

$$y_1^{(1)} = f\left(u_1, u_2, \mathbf{a}_{12}^{(1)}\right)$$

...

$$y_{m-1}^{(1)} = f\left(u_1, u_m, \mathbf{a}_{1m}^{(1)}\right)$$

$$y_m^{(1)} = f\left(u_2, u_3, \mathbf{a}_{23}^{(1)}\right) \tag{14}$$

...

$$y_{2m-3}^{(1)} = f\left(u_2, u_m, \mathbf{a}_{2m}^{(1)}\right)$$

...

$$y_s^{(1)} = f\left(u_{m-1}, u_m, \mathbf{a}_{m-1,m}^{(1)}\right)$$

The constant parameters of any function $f$ are optimized (separately for each neuron) by means of the LSM algorithm or the other selected training rule. Training process is always performed by using the subset of learning data.

The selection of the component elements for their processing accuracy is performed before the formed layer is connected to the network. The elements, which, on the basis of the selected evaluation criterion have too big processing error, are removed. $Q(y)$ (compare Fig. 4). There are various methods of performing the selection procedure:

*constant population method* is based on selection of $k$ neurons, for which $Q(y)$ reaches the least values,

*optimal population method* is based on rejecting these neurons for which the defined processing error is bigger than arbitrarily determined frontier $\square$ (usually $\square$ is determined separately for each layer on the basis of the least value $Q(y)$, the largest value $Q(y)$, or the range of values $Q(y)$ determined for a given generation of neurons),

*decreasing population method* defines the maximum number of elements in each layer and it will decrease along with the growth of the network.

Fig. 4. The illustration of the process of neuron selection in input layer


In the described way the neurons that model the searched input-output dependence most accurately are included in each layer, and the elements that introduce too big processing error are removed. The approved selection criterion $Q(y)$ plays a role of a mechanism of the structural optimalization at the stage of construing the new layer of neurons.

The definition of the evaluation criterion of the transition function is a very essential characteristic feature of the GMDH type of network. It allows any neuron to define the quantity of a processing error. In most criterion functions the division of the empirical data base into two independent parts is taken into account (Pham and the others 1995), one (learning data) is used by learning algorithm to optimize the constant parameters of the transition functions (see chapter 2.2), and the other (testing data) – with the evaluation criterion to calculate the introductory error. Thanks to this solution there is always certainty that it is internally not contradictory (Nagel and Newman, 1966). The risk of tuning in the network to the defined set of learning data instead of to the generalized rules represented by them is minimized this way.

Out of so many well known definitions of criterion functions (Farlow 1984, Ivakhnenko 1982, Ivakhnenko and the others 1987) some of them are intended for the general applications because of their features, the others – only for the selected applications). The most frequently used ones are:

1. *Regularity criterion*

$$\Delta^2(T) = \frac{\sum_{i=1}^{n_T}\left(y_i^* - y_i\right)^2}{\sum_{i=1}^{n_T} y_i^2} \tag{15}$$

where:
$n_T$ - amount of subset of testing data,
$y_i$ - measured values of output signal $y$,
$y_i^*$ - values of signal $y$ estimated by means of tested neuron.
Regularity criterion is a classical example of using the Godl's consistency criterion. The result of this application is the selection of neurons due to their best modeling of the trends of the signals' changes.

## 2. *The least deviation criterion*

$$n_{odch} = \frac{\sum\limits_{i=1}^{n}\left(y_{iL}^{*} - y_{iT}^{*}\right)^{2}}{\sum\limits_{i=1}^{n} y_{i}^{2}} \qquad (16)$$

where:
$n_L$ - amount of subset of learning data,
$n_T$ - amount of subset of testing data,
$n=n_L+n_T$ - amount of data set,
$y_i$ - measured values of output signal $y$,
$y_{iL}^{*}$ - values of signal $y$ estimated by means of learning neuron for subset $L$,
$y_{iT}^{*}$ - values of signal $y$ estimated by means of learning neuron for subset $T$.
The result of the neuron selection according to the least deviation criterion is accurate modeling of the value of signal $y$ for the given stimulation **u**. The characteristic feature of the solution is the synthesis of the testing neurons to verify the accuracy of the processing of the learning elements by means of the subset $L$. The drawback of the definition  (2.16) is the extended period of calculations connected with the construction of testing neurons and the impact of the result on the division of data set into parts $L$ and $T$.

## 3. *Convergence criterion*

$$i^{2}(n) = \frac{\sum\limits_{i=1}^{n}\left(y_{i}^{*} - y_{i}\right)^{2}}{\sum\limits_{i=1}^{n} y_{i}^{2}} \qquad (17)$$

where:
$n=n_L+n_T$ - amount of data set,
$y_i$ - measured values of output signal $y$,
$y_i^{*}$ - values of signal $y$ estimated by means of the testing neuron.
The dependence (2.17) is one of the few internal criteria not requiring the division of the empirical data set. The whole set $N$ is used for learning the neurons and testing their processing error. The convergence criterion is mainly applied in the interpolation tasks when the measurement results are not charged with big noises and measurement errors.
*Combined criterion*, such as

   *the least deviation plus convergence*

$$\rho_{1} = \sqrt{n_{odch}^{2} + i^{2}(n)} \qquad (18)$$

particularly useful in analysis of discrete systems

or

   *the least deviation plus regularity*
$$\rho_{2} = \sqrt{n_{odch}^{2} + \Delta^{2}(T)} \qquad (19)$$

55

often used in relation to the static cases.

The selected criterion (e.g. $Q(y)=n_{odch}$ , $Q(y)=i^2(n)$, $Q(y)=\square_{\square}$ or the others) is used to define the processing error of every (already trained) neuron.  On this basis the decision about its incorporation or rejection to the subsequent layer is made. Every new layer of the GMDH network  construed in a described way is not subject to any further changes.

In the second and all the next iteration of the GMDH procedure the output signals from the former layer are used as the input data. The neurons of the *l*-th (internal) layer are described by transition functions respectively which in the simplest case of only two input signals has a form of equations (20).

$$y_1^{(l)} = f\left(y_1^{(l-1)}, y_2^{(l-1)}, \mathbf{a}_{12}^{(l-1)}\right)$$
$$...$$
$$y_{s-1}^{(l)} = f\left(y_1^{(l-1)}, y_s^{(l-1)}, \mathbf{a}_{1s}^{(l-1)}\right)$$
$$y_s^{(l)} = f\left(y_2^{(l-1)}, y_3^{(l-1)}, \mathbf{a}_{23}^{(l-1)}\right) \qquad (20)$$
$$...$$
$$y_{2s-3}^{(l)} = f\left(y_2^{(l-1)}, y_s^{(l-1)}, \mathbf{a}_{2m}^{(l-1)}\right)$$
$$...$$
$$y_t^{(l)} = f\left(y_{s-1}^{(l-1)}, y_s^{(l-1)}, \mathbf{a}_{s-1,s}^{(l-1)}\right)$$

Analogical  to the dependence (12), the number of *k*-input neurons formed in the *l* layer  equals:

$$t = \binom{s}{k} = \frac{s!}{k! \left(s-k\right)!} \qquad (21)$$

where *s* is the number of elements selected in the layer *l*-1.

After forming the subsequent *l*-th layer of network the procedure of selection of the best (analogical as illustrated in Fig 2.4) is repeated.

On the same basis the neurons of further layers are formed until the so called *optimalization criterion* is met (Ivakhnenko 1987). The criterion is based on achieving in a construed layer of the network the minimum values of the selected criterion function $Q_{opt}$ . The least processing error $Q_{min}^{(l)}$ is calculated after completing every stage of the synthesis of the GMDH network (selecting s neurons and incorporating them into the network as the layer *l*)

$$Q_{min}^{(l)} = \min_{i=1,...,s} Q_i^{(l)}(y) \qquad (22)$$

then

$$Q_{opt} = \min_{j=1,...,l+1} Q_{min}^{(j)} \qquad (23)$$

Values $Q_i^{(l)}(y)$ can be defined by means of evaluation criterion chosen for selection purposes or the others. In the described way the evolution process of the network is run until the processing error starts to grow (See Fig. 5).

Fig. 5. The illustration of using the optimalization criterion of the GMDH network

In case of the combined criteria the achieved minimum is usually more „sharp". They also constitute one of the securities preventing the unconscious choice of the sub-optimal solution connected with the local minimum (see Fig.6a). In such cases or when the processing error only decreases asymptoticaly (this is often the case in using internal criteria or when noises and measuring errors influence, to some extent, the examined signal) the value of the threshold error for which the optimalization criterion is approved (see Fig. 6b) may be determined.



Fig.6. A typical cases of determining the optimalization criterion:  a) the use of the combined criterion to avoid sub-optimal solution,   b) the use of the value of threshold error with asymptotically decreasing criterion $Q_{min}$

The layer $l_{opt}$ , where the value of criterion $Q_{opt}$ has been achieved, plays a role of one element output layer. The output of the network is then a neuron for which $Q_i^{(l_{opt})}(y) = Q_{opt}$. The remaining elements of layer $l_{opt}$ are removed. As a consequence the neurons of the internal layers are also removed (from $l_{opt}$-1 to 2), that form blind branches not leading to the output cell

(see Fig. 7). The procedure of reduction of the elements not used to determine the output signals often considerably reduce the size of the network. At the same time it ends the process of synthesis.



Fig. 7. Final reduction of not used neuron of the GMDH type

The control of the process of the neuron network synthesis by calling off the external criteria is one of the most important assumptions of the Group Method of Data Handling.

In practice the approximation of consistency theory is applied (Nagel and Newman, 1966), where the set of the empirical data is divided in two parts: learning and testing. The third subset of *controlling data* used for better tuning in of the final solution is distinguished in some cases where the accuracy of the network model is of more importance than the time of calculations, and the set of the measuring samples is considerably large.

According to the Gödl's theory the controlling data should be used as the basis of determining the optimalization criterion independent of the criteria used in the former stages of the network synthesis. Then, first subset of data (learning) is used for learning the particular neurons, second (testing data) – for selection of the elements in a newly construed layer, and the third (controlling data) – for evaluation of the effectiveness of processing the whole formed structure. Such a solution assures the optimal compromise between the accuracy of the estimation of the measuring samples and the features of the generalization of the neuron network.

The other method of using the controlling data prepared by Adaptronics Corporation (Hecht-Nielsen 1991) is based on setting the equations of regression for Ivakhnenko's polynomial, modeled by the construed GMDH network. Then the correction of the constant parameters is possible by means of the least mean squares' method. As the presented technique is based on tuning in all the parameters at the same time, the proper amount of regression equations is required for its performance i.e. the proper large subset of the controlling data. The practical applications confirm that the presented method of extra learning of the ready network increases the accuracy of the estimation of the output signal *y*.

## 3. Generalizations of the GMDH algorithm

The assumptions of the neuron networks of the GMDH type presented in chapter 2 give lots of freedom in defining the particular elements of the algorithm of synthesis. The mentioned possibilities relate to, for example, the definition of the transition function, learning rules, evaluation criterion of the processing accuracy, optimalization criterion, etc. The concept of the group method of data handling also allows developing the formula of the GMDH network due to the supplement with the additional procedures that improve the accuracy of modeling and extending the scope of applications.

Some of the mentioned modifications of the algorithm of the network synthesis are based on using an additional (third) subset of the controlling data. The main idea is an increase of the accuracy of the output signal estimation and also the improvement of the features of the generalization information hidden in the set of the empirical data. The optimized resulting structure of the network allows using it in a more extensive analysis of the examined object. The concept of finding the optimal structural and parametrical solution by means of the evolution of the simple component elements (neurons) is maintained here.

Freedom of defining certain elements of the procedure of the GMDH network construction also includes the generalizations that extend the scope of their applications. The scheme of the algorithm of synthesis allows extending its formula to achieve new quality of solutions while maintaining the advantages resulting from the structural and parametrical optimalization of the GMDH.

## 3.1 System of many inputs and many outputs

The classical assumptions of the group method of data handling presented in chapter 2 always lead to the formation of the neuron network of many inputs ($u_1$, $u_2$, ..., $u_m$) and one output ($y$). The construed structure approximates the dependence (24).

$$y = f\left(u_1, u_2, ..., u_m\right) \tag{24}$$

The systems of many inputs ($u_1$, $u_2$, ..., $u_m$) and many outputs ($y_1$, $y_2$, ..., $y_p$) are found in the practical applications most often. In this case any output quantity is modeled by means of a separate GMDH network, even if all the examined stimulation signals and reactions are connected with each other. In the classical solutions of the neuron networks (Korbicz and others 1994, Pham and others 1995) the resulting topology has a form of many inputs – many outputs, that model the whole examined system instead of one of the described dependencies only. The introduction of a small correction in the procedure of the neuron selection can lead to the identical solution also in case of the network of the GMDH type.

The mentioned modification relates to the method of computation of the processing error. There are two methods to solve this problem that depend on the synthesis and the size of the resulting structure as well as processing accuracy. The use of the network of the GMDH type in the parallel estimation of many parameters ($y_1$, $y_2$, ..., $y_p$) creates a necessity to calculate more then one criterion of the quality control. The value of the criterion of the estimation accuracy $Q^{(k)}(y_1)$, $Q^{(k)}(y_2)$, ..., $Q^{(k)}(y_p)$ (calculated for any output quantity - see Fig. 3.1) can be determined for any neuron generated in layer $k$.

Fig. 8. A scheme of determining the processing error in a newly formed neuron layer

In a simple solution the values $Q^{(k)}(y_1)$, $Q^{(k)}(y_2)$, ..., $Q^{(k)}(y_p)$ are used to determine the generalized processing error $Q^{(k)}$. It can be determined for any construed neuron $Q^{(k)}$ by one of the methods presented in table 1.

Tab. 1. Selected methods to determine the neuron's processing error in a multi-output GMDH network

| Method | Calculation formula | |
|---|---|---|
| Minimum error | $Q^{(k)} = \min_{i=1,...,p} Q^{(k)}(y_i)$ | (25) |
| Maximum error | $Q^{(k)} = \max_{i=1,...,p} Q^{(k)}(y_i)$ | (26) |
| Medium error | $Q^{(k)} = \dfrac{1}{p} \cdot \sum_{i=1}^{p} Q^{(k)}(y_i)$ | (27) |
| Medium square error | $Q^{(k)} = \sqrt{\dfrac{\sum_{i=1}^{p} \left(Q^{(k)}(y_i)\right)^2}{p}}$ | (28) |

The effectiveness of the activity of every neuron can be evaluated only by one criterion $Q^{(k)}$ in this way. This allows performing the selection of neurons in the particular layers according to the scheme described in chapter 2. The optimization criterion is also used in accordance with the definition (23). The partial components of the processing error $Q^{(k)}(y_1)$, $Q^{(k)}(y_2)$, ..., $Q^{(k)}(y_p)$ are used after the process of synthesis is completed to perform the classification of outputs in the last layer $k=l_{opt}$ (see Fig. 9). The particular output signals $y_i$ are matched with the output of the neuron for which $Q^{(l_{opt})}(y_i)$ reaches the least value.

Fig. 9. Scheme of the classification of the output neuron of the GMDH network

The elements of the last layer that are not used, are removed. As a result the neurons of the internal layers that do not take part in the generation of any of the output signals are also reduced. After performing the described procedure the network has a structure as presented in Fig. 10. It is characterized by a compact topology with all the outputs placed in the same layer.



Fig. 10. The multi-output network of the GMDH type structure

In the above-described solution the specialization of the neurons to estimate the concrete signals is performed in the last layer. The solution performing this procedure from the beginning of the process of the network synthesis is possible. Then, the independent evaluation of any of the criteria $Q^{(k)}(y_1)$, $Q^{(k)}(y_2)$, ..., $Q^{(k)}(y_p)$ is performed after the generation of each layer of the neurons (see Fig. 3.4). The elements that introduce too big estimation error of each output $y_1$, $y_2$, ..., $y_p$ are removed. The effectiveness of the neuron in processing at least one output signal is sufficient to leave the neuron in the network.

$y_1^{(k-1)} \longrightarrow f \longrightarrow y_1^{(k)}$ : $Q_1^{(k)}(y_1) > \varepsilon$, $Q_1^{(k)}(y_2) < \varepsilon$, ..., $Q_1^{(k)}(y_p) > \varepsilon$

$y_2^{(k-1)} \longrightarrow f \longrightarrow y_2^{(k)}$ : $Q_2^{(k)}(y_1) > \varepsilon$, $Q_2^{(k)}(y_2) > \varepsilon$, ..., $Q_2^{(k)}(y_p) < \varepsilon$

$y_3^{(k-1)} \longrightarrow f \longrightarrow y_3^{(k)}$ : $Q_3^{(k)}(y_1) > \varepsilon$, $Q_3^{(k)}(y_2) > \varepsilon$, ..., $Q_3^{(k)}(y_p) > \varepsilon$ $\Rightarrow$ *neuron removed from layer (k)*

$y_4^{(k-1)} \longrightarrow f \longrightarrow y_4^{(k)}$ : $Q_4^{(k)}(y_1) > \varepsilon$, $Q_4^{(k)}(y_2) > \varepsilon$, ..., $Q_4^{(k)}(y_p) > \varepsilon$ $\Rightarrow$ *neuron removed from layer (k)*

$y_l^{(k-1)} \longrightarrow f \longrightarrow y_s^{(k)}$ : $Q_s^{(k)}(y_1) < \varepsilon$, $Q_s^{(k)}(y_2) > \varepsilon$, ..., $Q_s^{(k)}(y_p) > \varepsilon$

Fig. 11. Synthesis of the neuron layer in the multi-criteria selection

The optimalization criterion in the proposed solution is applied independently to any value of the partial error ($Q^{(k)}(y_1)$, $Q^{(k)}(y_2)$, ..., $Q^{(k)}(y_p)$). The synthesis process is complete once any of the calculated criteria values reaches minimum. Similarly to the case described earlier, the output signal $y_1$ is connected with the output of this neuron for which $Q^{(k)}(y_1)$ achieved the least value, $y_2$ – with the output of this neuron for which $Q^{(k)}(y_2)$ achieved the least value etc. The difference is that the particular minimum could occur in different stages of the network synthesis. This is why in the described solution the outputs of the resulting structure are usually in different layers (see Fig. 12).



Fig. 12. The structure of the GMDH multi-output network constructed by the method of the multi criteria selection

## 3.2 Dynamic systems

The procedure of the synthesis of the artificial neuron network of the GMDH type was created on the basis of the algorithm of the polynomial modeling. This is why its main application is the estimation the unknown output signals **y** on the basis of input signals **u** measured in a definite

moment of time $t_k$. Taking into account that the learning process is based on using the collected measuring samples (stimulation $x_i$ and reactions $y_i$ measured in moments $t_i < t_k$), the network model of the GMDH type can be interpreted as approximation of dependencies (29).

$$y_{k+1} = f\left(y_k, y_{k-1}, y_{k-2}, ..., u_{k+1}, u_k, u_{k-1}, u_{k-2}, ...\right) \tag{29}$$

where:

$u_k$, $y_k$     - current values of signals $u$ and $y$,

$u_{k+1}$, $y_{k+1}$ - next (future) values of signals $u$ and $y$,

$u_{k-i}$, $y_{k-i}$    - former values of signals $u$ and $y$.

The neuron network modeling the dependence (29) can have an application in the analysis of the static systems. In relation to the dynamic systems the scope of applications is limited to the controlling systems in which the stimulation values in the next moment of time $t_{k+1}$ are well known. The presented analytical apparatus can also be used to calculate the unknown value of the signal $y_{k+1}$ at the moment $t_{k+1}$ on the basis of the signals $y_k$, $y_{k-1}$, $y_{k-2}$, ... and $u_k$, $u_{k-1}$, $u_{k-2}$, ... measured in moments $t_k$, $t_{k-1}$, $t_{k-2}$, ...,thus – to estimate the dynamic systems. In this case the GMDH network is used to solve the prediction tasks (30).

$$y_{k+1} = f\left(y_k, y_{k-1}, y_{k-2}, ..., u_k, u_{k-1}, u_{k-2}, ...\right) \tag{30}$$

For the synthesis of the neuron network of the GMDH type solving the problem specified by the dependence (30) the generation process of the particular neuron population and their selection should be held identically as described earlier. The only modification necessary in case of the other type of the processing data relates to the definition of the GMDH neuron.

The learning and testing data are measured in the discrete moments of time. As a consequence – the estimated signals are also described in the discrete moments. The prediction tasks of the discrete dynamic systems in the network model of the GMDH type are most often solved due to the evolution of the neuron layers defined by the transition function determining a certain template of dependencies (time frame) between the input signals. In case of one stimulation $u$ the transition function has a form (31).

$$y_{k+1} = f\left(u_k, u_{k-i}, u_{k-j}, ...\right) \tag{31}$$

In a general case, the estimated signal $y_{k+1}$ can depend on the former values of more than one stimulation. Therefore the time frame modeled by one single neuron should include different input signals $u_1$, $u_2$, ..., $u_s$ registered in different discrete moments of time $t_k$, $t_{k-1}$, $t_{k-2}$, ..., $t_1$. The structure of the GMDH neuron intended for modeling of the dynamic systems has a form as presented in Fig. 13.



Fig. 13. The structure of the GMDH neuron defined for the dynamic discrete systems

The transition function of the neuron described in Fig. 13 must meet the condition (2) and the linearity condition towards the constant parameters as it happens in any variant of the GMDH type. Therefore the form of the function $f$ should be in accordance with the definition (32).

$$y_{k+1} = a_0 + a_1 \cdot f\left(u_{1,k-i}\right) + ... + a_2 \cdot f\left(u_{2,k-j}\right) + ... + a_s \cdot f\left(u_{s,k-l}\right) + ...$$
$$+ a_{12} \cdot f\left(u_{1,k-i}, u_{2,k-j}\right) + ... + a_{1s} \cdot f\left(u_{1,k-i}, u_{s,k-l}\right) + ... + a_{2s} \cdot f\left(u_{2,k-j}, u_{s,k-l}\right) + ... \qquad (32)$$

From the definition of the GMDH neuron determined for the discrete dynamic systems results the necessity to apply the recurrent procedure with the estimation of the signals $y_{k+2}$, $y_{k+3}$ etc. In these situations the modification of the transition function that takes into account the time as an additional processing signal is sometimes applied. Then, the signal $y_{k+i}$ is estimated in one course of the group method of the data handling (without the recurrent procedure) for any given time $t_{k+i}$. This solution brings better results also in the situations where the sampling period is non-uniform or when the measuring samples for some moments of time are not registered in database. The mentioned alternative concept of the dynamic GMDH neuron leads to the structure presented in Fig. 14.



Fig. 14. The structure of the dynamic GMDH neuron synchronized by the parameter of time

The neuron transition function presented in Fig. 14 is determined by the dependence (33).
$$y = f\left(u_1, u_2, ..., u_s, t\right) \qquad (33)$$
It has a polynomial form in a classical approach (33). Therefore in case of the two input signals $u_1$ and $u_2$ (three-input dynamic neuron) according to the definition (4) $f$ can be described as the polynomial of a second degree.
$$y = a_0 + a_1 \cdot u_1 + a_2 \cdot u_2 + a_3 \cdot t + a_{11} \cdot u_1^2 + a_{22} \cdot u_2^2 + a_{33} \cdot t +$$
$$a_{12} \cdot u_1 \cdot u_2 + a_{13} \cdot u_1 \cdot t + a_{23} \cdot u_2 \cdot t \qquad (34)$$
It is worth pointing out that if $t = const.$, then the dependence (34) is identical with the function (4). Therefore in the analysis of the constant processes, the modeling of the static systems by means of the neuron network of the GMDH type is a special case of the general procedure of the network synthesis modeling the dynamic systems. The area of using the presented analytical apparatus is only dependent on the way the time is determined: as variable or constant.

The two presented variants of the dynamic neuron network of the GMDH type differentiate due to the way of the synchronization of the processing input signals in relation to time. In one case it is a defined time frame, in the other – time as a directly processed parameter. In both cases the GMDH network is constructed from the multi-input neurons. This is connected with the extended learning time in comparison with the static case. If the least mean squares method is used to train neurons it is necessary to perform multi-measured matrix calculation many times. In case of using one of the classical learning algorithms it is necessary to have the set of empirical data in a right amount and the arduous calculation procedure.

The greater scope of the applications of the neural network of the GMDH type than modeling only is possible in relation to the dynamic systems. The result of using the polynomial transition function is the construed networking model that allows the identification of the function (Ivakhnenko polynomial) modeling the dynamics of the examined object. The characteristic polynomial is easy to get after using the transformation Z (35).

$$a_n \cdot z^n + \ldots + a_2 \cdot z^2 + a_1 \cdot z + a_0 = 0 \tag{35}$$

The dependence (35) allows using the criteria of stability examination of the modeled dynamic system known from the controlling theory (Ivakhnenko 1971). In accordance with these criteria the roots of an equation (35) placed inside the circle of an elementary radius on a plane of complex variables is the condition of stability.

## 3.3   Harmonic systems

The identification of the modeling function by the GMDH network enables to perform the spectral analysis of the examined system in a similar way (Ivakhnenko and others 1987) with the transition functions defined by means of sinusoidal dependence on time e.g. (36) or (37).

$$y = f_h(u) = a_0 + a_1 \cdot u \cdot \cos(\omega_1 \cdot t + b_1) \tag{36}$$

$$y = f_h(u) = u \cdot (a_0 \cdot \cos \omega \cdot t + a_1 \cdot \sin \omega \cdot t) \tag{37}$$

According to the presented structure, the frequency analysis in a described class of network is always connected with the dynamics of signals. The constant sampling period is an essential condition (in contrast to the classical polynomial approach).

The structures formed from neurons processing the signals in accordance with the dependence of form (36) constitute a separate class of the GMDH network of the specific characteristics that are applied particularly in the examination of the frequency parameters of the dynamic objects. These types of solutions are named *GMDH harmonic network* due to the variation of the characteristic features.

In the synthesis process the neuron networks of the harmonic type are, similarly to the classical solution, constructed each separately and are incorporated into the neuron structure in the groups forming the layers. Determining the constant parameters of every neuron is the key problem for realization of this task.  It is not possible to adapt the methods used in the polynomial transition function because the linearity condition towards the constant parameters is not met.

The approximate algorithm to calculate the parameters of the models is well known only. (3) (Ivakhnenko and others 1975). According to its assumptions it is necessary to determine the number *m* of neurons in a formed layer following the known criteria (Banaszak and Kuś 1993, Farlow 1984). Then it is possible to demonstrate that for the parameter $\square_j$ minimizing the medium square criterion (38)

$$\Delta^2 = \sum_{k=m+1}^{N_L - m} \left[ \sum_{j=1}^{m} \alpha_j \cdot (y_{k+j} + y_{k-j}) - y_k \right]^2 \tag{38}$$

where:  $N_L$ – size of learning data subset
is formed and equation (39)

$$2 \cdot \sum_{j=1}^{m} \alpha_j \cdot \cos \omega j = 1 \tag{39}$$

that by means of the formula (40) (Bronstein and others 1995)

$$\cos \omega j = 2 \cdot \cos \omega (j-1) \cdot \cos \omega - \cos(j-2)\omega \tag{40}$$

can be transformed into a form of the polynomial equation of *m* degree towards the variable $\cos\square$. The solutions of this equation enable to determine parameter $\square_j$ for every neuron of the construed layer. The parameters $a_0$ and $a_1$ are determined, alike in a GMDH classical approach (Hecht-Nielsen 1991) because they meet the linearity condition.

The process of forming and connecting the neuron layers in the harmonic network differentiates from the well known classic approach. The first difference is the assumption of a number of the elementary cells *m* generated in an input layer. The value *m* is dependent on the assumed

accuracy of the spectrum analysis. The other characteristic feature of the described type of the network is the usage of two types of neurons: harmonic and summing one. The means of construing harmonic neuron has been described in section 3. Therefore, the summing neurons are described by the transition function $f_s$ (41).

$$y = f_s(\mathbf{u}) = \sum_{i=1}^{n} a_i \cdot u_i \qquad (41)$$

where: $n$ – a number of neuron inputs.

In a simple case $n = 2$. Constant parameters $a_i = 1$ are often taken into consideration because amplification coefficient is also included in components $a_0$ and $a_1$ of the harmonic neurons.
Layers of the GMDH network are formed alternately using first only harmonic cells and next - only summing cells. After constructing the harmonic layer, the summing neurons are formed for all the possible combinations of the output signals. The coefficient of transition accuracy $Q = \square^{\square}$ in accordance with accepted criterion (Farlow 1984) is calculated for each summing neuron. Then the computed values of error transition are used to perform the neurons selection (a summing layer). The described process of synthesis of neuron layers is illustrated in Fig. 15.



Fig. 15. Synthesis process and selection of the neurons in the GMDH harmonic network

In the described procedure, neurons, which in the most precise way model the examined input-output dependence, are incorporated to each summing layer, whereas those elements which introduce too much of a transition error are removed. In this way these layers also play the role of a filter for the harmonic components, which contain too much of an error. The accepted criterion of selection plays the role of mechanism of structural optimization at a stage of construing the new neuron layers.

## 4.  Network applications of the GMDH type

The assumptions of the artificial neural network of the GMDH type presented in chapter 2 and 3 show that the group method of data handling due to a considerable freedom to define some elements of the synthesis algorithm offers a wide scope of possible applications. The use of the classical solutions of the described algorithm gives an opportunity to estimate the static systems and automatic steering systems. Due to the introduction of some generalizations in the neuron definition and the procedure of their selection the practical applications of the GMDH network

can also be applied in relation to the complex dynamic systems. The characteristic features for the presented classes of the neuron network are easy to be identified in these applications.


## 4.1    Damage classification in the electronic system

The features of the group method of data handling such as small requirements for the size of the learning data set, short learning time, high accuracy of estimation due to the structural and parametrical optimalization of the neuron network allow the described analytical apparatus to be applied also in the systems of real time. Such systems are among others the diagnostic systems using the principles of analytical redundancy. The comparison of the estimated signal with the measured value enables to identify the state of the examined object as normal or emergency one. The attempts to use the neuron network of the GMDH type to detect the damages were performed in the analysis of the feeder function of the digital and analogue electronic systems. 35 samples of voltages and currents measured in the different points of the system and in the different state of its normal labor were recorded in the database. The chart presented in Fig 16 illustrates the course of variability of a given residuum (the difference between the measured and estimated signal) in the experiment, where the damage of the output circuit (caused by the operator) –15V took place.



Fig. 16. The course of the residuum variability for the output circuit –15V electronic feeder


The received results show that the collected empirical data base turned out to be sufficient to construct the network model accurate for the estimation of the signal $V_{-15V}$ and consequently for the precise classification of the damage state.

The qualities connected with the use of the empirical data and inductive generalization of the hidden information about the examined regime are more visible in the applications of the GMDH network than in the other neuron networks. These features have been tested in the performed diagnostic experiments by means of artificial introduction of the measuring noises. The results presented in Fig. 4.2 show that the connection of the mentioned features of the network synthesis algorithm based on the group method of data handling ensures high resistance of the resulting structure to the disturbances of the measured signals.

Fig. 17. The charts of variability of the selected residues and their dependence on the level of the measuring noises

The results presented in Fig. 17 give evidence that the mechanism of the GMDH structural and parametrical optimalization supported by the concept of the internal consistency of the solution (by means of using the separate subsets of learning and testing data) leads to more effective compromise between the accurate approximation of the signals value and correct modeling of the trends and their changes.

## 4.2 Modeling selected processes of the sugar industry

The classical form of the algorithm of the GMDH neuron network synthesis is useful in modeling tasks of the static systems or control systems (see dependence (29)) The performed experiments show that the extension of the scope of the possible applications while maintaining the features connected with the accurate modeling based on small sets of learning data is the result of generalization and modifications of this theory that do not break the basic scheme. In particular the mentioned generalizations can influence the analysis of the dynamic systems.

The experiments were related to the examination of the different elements of the „Lublin" technological sugar factory course (boiler house and evaporating station). The basic goal was to control the accuracy of the dynamics modeling of the considered processes. The multi-elements set of the learning data was used in the synthesis of the modeling network. The GMDH network was used for prediction of the selected quantities in advance by one measuring step during the trials. The estimated values were compared with the measured ones. The exemplary results are illustrated in Fig. 18.

68

Fig. 18. The course of the variability of the measured juice level and adequate estimated values
The experiments run for the sugar processes confirm the capabilities of the GMDH network to model accurately both the values of the measured signals and trends of their changes in case of the dynamic systems also.


## 5.  Conclusions

The outlined theory of the network of the GMDH type enables to solve some problems that cause the restrictions of the classic solutions of the artificial neural networks. The features opening the new area of applications and introducing new quality of the neuron calculation theory were achieved due to the integration of the other known methods used in the modeling tasks.

The elements of the arbitrary structure definition (topology), which is a feature (and at the same time one of the weaknesses) in the classical approach, are missing in the presented solution. The problem of determining the network topology as an additional learning effect has been solved. What's more – the group method of data handling always leads to the optimal structure for the considered problem. The GMDH network grows and evolves in the learning process until development leads to the improvement of the effectiveness of activity.  Taking into account the fact that the means of optimalization of the parameters of the constant transition functions (e.g. LSM method) was defined, the proposed method leads to the structural and parametrical optimalization of the neuron networks.

The inconsistency of the traditional solutions based on using the deductive methods (to determine the number of neurons and the internal layers) in the mathematical apparatus based on inductive conclusion has been eliminated in the presented way. The elements of the Gödl's inconsistency theory that fulfilled the mathematical apparatus show few examples of its practical application and considerably improves the qualities of generalization of the features hidden in the empirical data of the examined object. The analytical mechanism construed this way merges both quantity and quality approach in the modeling tasks.

The essential quality of the presented method of the neural network synthesis is not only capability to accurate estimation of the values of the analyzed signals and trends of their changes but also to identify the structure of the mathematical model. The variant of the harmonic GMDH network is an example of using the mentioned feature. Due to the specific construction of the formed mathematical model it is possible to infer about the physical qualities of the examined object, for example, about its stability and – in case of the harmonic network – about its frequency parameters.

Paradoxically the elements that decide about higher effectiveness of the GMDH type of network cause the restrictions of the scope of the practical applications in comparison with the classical solutions. The synthesis of the neural network using the group method of data handling requires

69

measuring data base and expertise knowledge that is important to make decisions e.g. about the form of the transition function of the generated neurons (polynomial, sinusoidal etc) of the criteria functions etc.

Taking into account the extensive possibility of selection of the transition functions, the criteria functions, methods of neuron selection, algorithm of the division of the learning and testing data, one can say that the GMDH networks constitute a separate class of solutions with the interesting features. The new qualities allowing the extension of the scope of applications also in the built-up dynamic systems, both linear and non-linear are the result of the applications of the structural and parametrical procedure of the optimalization of the network. The presented solution is efficient in the hardware applications, computer implementation and applications of the real time due to the effective processing of small sets of learning data and short time of synthesis even of the multi-layer networks. The presented qualities and considerable freedom to define some elements of the algorithm show perspective of the development of the theory of the neuron networks the GMDH type.

## References

1.  Anderson, B.D.O., Moore, J.B., 1984, Filtracja Optymalna, WNT, Warszawa.

2.  Banaszak, Z., Kuś, J., 1993, Group method of data handling in technical diagnostic tasks, *Applied Mathematics and Computer Science*, Vol. 3, No. 3, pp. 573-593.

3.  Bronstein, I.N., Semendiaev, K.A., 1995, Matematyka, PWN, Warszawa.

4.  Duran, B.S., Odell, P.L., 1974, Cluster Analysis, Springer Verlag, Berlin.

5.  Farlow, S.J. (Ed.), 1984, Self-organizing Methods in Modelling: GMDH-type Algorithms, Statistics: Textbooks and Monographs, V. 54, Marcel Dekker Inc., New York.

6.  Hecht-Nielsen, R., 1991, Neurocomputing, Adison-Wesley Publishing Co., New York.

7.  Ivakhnenko, A.G., 1971, Polynominal Theory of Complex Systems, *IEEE Trans. on Systems, Man and Cybernetics,* Vol. SMC-1, No. 4.

8. Ivakhnenko, A.G., 1975, Dolgosrochnoje prognozirowanie i uprawlenie sloznymi sistemami, Technika, Kijów.

9. Ivakhnenko, A.G., 1982, Induktiwnyj Metod Samoorganizacji Modelej Sloznych System, Naukowa Dumka, Kijów.

10. Ivakhnenko, A.G., Jurackovskij, J.P., 1987, Modelirowanie sloznych system po eksperimentalnym dannym, Radio i Swiaz, Moskwa.

11. Kohonen, T., 1984, Self-organization and Associative Memory, Springer Verlag, Berlin.

12. Korbicz, J., Obuchowicz, A., Uciński, D., 1994, Sztuczne Sieci Neuronowe, Akademicka Oficyna Wydawnicza PLJ, Warszawa.

13. Nagel, E., Newman, J.R., 1966, Twierdzenie Gödla, Współczesna Biblioteka Naukowa OMEGA nr 52, PWN, Warszawa.

14. Pham, D.T., Xing, L., 1995, Neural Networks for Identification, Prediction and Control, Springer Verlag, London.

15. Tadeusiewicz, R., 1993, Sieci Neuronowe, Akademicka Oficyna Wydawnicza RM, Warszawa.

# Medium structure modeling on parallel computers[1]

*Marcin Paprzycki, Valerii Rozenberg and Boris Digas*

## Abstract

The work deals with application of methods of computer diagnostics to different problems related to assessment of natural and human-caused risk. Some aspects of usage of computational methods in the fields of plasma physics and geophysics are considered. A solution algorithm for the problem of a medium structure reconstruction is suggested. Some results of computational experiments are discussed.

## Introduction

Fast development of computer capabilities allowed for the development of *computer diagnostics* as a branch of computer modeling. Here, computer modeling can be considered to be a process in which the set of parameters of a model is adjusted to obtain a certain *ideal* result. Computer diagnostics, on the other hand, encompasses techniques for investigation of particular characteristics of an object, where the data used as input is a result of a measurement (and thus contains errors). One of the important features of the computer diagnostics is that it involves calculations with large amount of data and thus requires application of modern methods of computational mathematics and advanced computer hardware and software.

Application of computer diagnostics in medicine permits not only to increase the accuracy of the result but also facilitates and accelerates development of new methods of treatment. Usage of computer diagnostics for testing industrial output and machinery, industrial process and its control opens up possibilities of development of new technologies in the industry. Computer modeling is of great importance in environmental monitoring, astrophysics, geophysics, and various fields of physics, chemistry, biology etc.

Computer diagnostics may also serve as a background for different methodologies of catastrophic risk management. Both natural and human-caused catastrophic risk may be investigated by use of mathematical models and computer simulations. Presently, the advanced computational approaches allow us to deal with large-scale decision-making problems in the presence of multidimensional mutually dependent random variables [9].

As an example of dealing with the case of natural risk assessment one can consider application of computer diagnostics is support of data-based modeling approach to the decision-making in case of property insurance in seismically active regions. In this case decisions have to be made in a situation involving a substantial lack of data. More precisely, the existing historical data about the earthquakes is insufficient for predicting seismic events at any particular location, although rich data about their occurrence may exist on an aggregated regional level. Catastrophe modeling is aimed at compensation, to some extent, for the lack of historical data. Models have to play a key role in generating data used in designing new policies. Different catastrophes may exhibit a wide spectrum of impacts on the public health, environment and economy. Each of these episodes seems to be characterized by an extremely low probability and may be simply ignored in the so-called "practical approach".

Another methodological challenge involved in catastrophic risks management is related to the fact that catastrophes have different spatial patterns and quite differently affect various locations. The location of a property (private or industrial) with regard to the center of an earthquake is extremely important information. Together with the regional geology and the soil conditions the location influences the degree of vibrations and, hence, the damage incurred at a given location. Management of complex interdependencies among catastrophic risks, losses and decisions is possible only within a geographically explicit framework and it is the role of computer diagnostics to provide the necessary data for the decision-makers.

The need for innovative R&D involving computer diagnostics has been recognized by the nuclear industry and by those countries that believe in the overall benefits, viability and importance of nuclear power for the long-term solution of their energy problems [19, 20]. As it will be shown in this paper, there exists a number of similarities between the research involved in both nuclear physics and risk management when the problem of medium structure is considered. More precisely, in the present work, we will consider examples of application of computer diagnostics to the plasma physics and geophysics. Section 1 outlines main principles of diagnostics of high-temperature plasma. In Section 2, the problem of reconstruction of a medium structure is formulated and formalized. Section 3 describes a solution algorithm for this problem. Some results of computational experiments are discussed in Section 4.

# 1 Diagnostics of high-temperature plasma: problem statement

As indicated above we use the term *computer diagnostics* to describe the process of determining of a quantitative (or qualitative) characteristics *x* of a given object based on the results of measurements. These measurements are substituted for the "true" information resulting in the general transformation in the form *y = Ax*.



Fig. 1: Scheme of diagnostics of high-temperature plasma [17].

Let us consider an example. Different physical experimental methods of molecular spectroscopy, nuclear physics, and astrophysics have found their use in modeling of high-temperature plasma [22]. The substituted (measurement based) information y of the physical characteristics *x* of plasma is used. The aim of the diagnostics process is to find the correct value of the characteristics *x*. As an example of such a process consider X-ray based methods of determining of the spatial distribution of electron temperature of the plasma. Similarly, one can consider determination of the spatial distribution of electron density of the plasma by the method of interferometry or by the use of laser dispersion and other methods for the determination of particular parameters of the plasma (see, Fig. 1).

Let $x = f(r)$ be a spatial distribution of plasma's parameter under consideration. Let the plasma formation be cylindrically symmetrical. Let a detector register integral radiation from different elements of the plasma volume observed at an angle $\theta \in [0,\pi)$. Denote by $\varphi(x')$ the registered signal at $\theta \in [0,\pi)$. Axis $Ox'$ is directed at an angle $\theta$ to the axis $Ox$ and is connected with the direction of the signal registration. In this case, the relationship between the unknown distribution $f(r)$ and the signal $\varphi(x')$ has the form:

$$\varphi(x') = 2\int_{x'}^{R} \frac{f(r)r\,dr}{\sqrt{r^2 - (x')^2}}.$$

This relationship is the Abel's integral equation with respect to $f(r)$. So, even in the simple case of modeling of the cylindrically symmetrical plasma, one has to solve this equation along with performing physical measurements. Numerical solution of an equation of this type for inaccurate data is a rather complicated problem. A method for solving a problem that leads to a similar mathematical formulation will be described in detail below.

## 2 Reconstruction of a velocity field: problem statement

Let domain $\Omega \in \mathbf{R}^2$ be a model of a region on the Earth's surface containing a set of sources and receivers of seismic signals. Assume that the number of sources and receivers is *large enough*. Let $t_1,\ldots, t_n$ represent measured times of seismic signal propagation corresponding to the different pairs "source-receiver". It is assumed that the process of propagation of the signal is described by the laws of geometrical optics. In this case each value of the travel time of the signal is correlated with the velocity characteristics of the medium by the following integral relationship:

$$t_i = \int_{l_i} \frac{dl}{v(r)}, \quad (i = 1,2,\ldots,n),$$

where $v(r)$ is the velocity of propagation of seismic waves ($r = (x,y) \in \Omega$); $l_i$ is the integration contour corresponding to the $i$-th seismic ray. In this case our aim is to determine values of the function $v(r)$ in every point of the region under consideration.

This problem is nonlinear, but it can be linearized [1, 4, 18] if one considers the non-dimensional value

$$m(r) = \frac{v^{-1}(r) - v_0^{-1}(r)}{v_0^{-1}(r)}$$

as the function subject to the reconstruction. Here, $v_0(r)$ is a zero approximation of the seismic wave propagation velocity. Besides, one should substitute the integration contours with some approximations $l_{0i}$ of the unknown trajectories. In this case, one can determine the value $m(r)$ from the following relationships:

$$\int_{l_{0i}} \frac{m(r)}{v_0(r)} dl = \delta t_i \quad (i = 1,2, \ldots, N), \tag{1}$$

where $\delta t_i = t_i - t_{0i}$. Symbols $t_{0i}$ and $l_{0i}$ denote the travel times and the trajectories of the signals, calculated for the initial approximation $v_0(r)$.

Taking into account equalities (1), one can formulate the following mathematical problem for finding the distribution $m(r)$, $r \in \Omega$, which permits to pass easily to the velocity distribution under search. Given are values $\gamma_i$, ($i = 1,2,\ldots,n$), which can be represented in form of linear functionals of unknown function $m(r)$:

$$\gamma_i = \iint_{\Omega} G_i(r)m(r)\,dr, \tag{2}$$

where $G_i(r)$ are the kernels, whose form is determined by the initial data (shape of rays).

Quite often [1, 4, 11, 18], among all solutions satisfying (2), one selects a function that minimizes the quadratic functional

$$L(m) = \iint\limits_{\Omega} |m(r)|^2 \, dr \, . \tag{3}$$

Thus, the problem of velocity characteristics reconstruction may be formalized in the form of the problem of finding the function $\hat{m}(r)$ satisfying the constraints (2) in form of equalities and minimizing functional $L(m)$.

Note that, actually, travel times are measured inaccurately. Therefore, it is natural to consider the problem of minimization of the functional (2) with the constraints in form of the inequalities:

$$\delta t_i - \varepsilon^{(2)} \le \iint\limits_{\Omega} G_i(r)m(r) \, dr \le \delta t_i + \varepsilon^{(1)}, \quad i \in [1:n],$$

where $\varepsilon^{(1)}$, $\varepsilon^{(2)}$ are the measurement errors.

Denote by $(\cdot,\cdot)_\Omega$ and $\|\cdot\|_\Omega$ the scalar product and the norm of space $L_2(\Omega)$ respectively. Then, the above problem may be rewritten as the problem of finding a minimal $\|\cdot\|_\Omega$-norm solution of a system of linear inequalities in a Hilbert space:

$$\begin{cases} \min \|m\|_\Omega^2 \\ \delta t_i - \varepsilon^{(2)} \le (G_i, m)_\Omega \le \delta t_i + \varepsilon^{(1)}, \quad (i = 1,...,n). \end{cases} \tag{4}$$

The theory of linear inequalities is a very efficient tool for solving a wide range of problems [2, 5, 10, 23]. One of methods of solving problem (4) was suggested in [7,13]. This method, known as the constraint aggregation method, was later modified for a number of different situations (see, [3, 6, 8, 12, 14, 15, 16]).

In the present work, we will describe a new algorithm for solving the velocity reconstruction problem, which is based on the ideas of [7,13]. According to those, we solve the extremal problem by use of an iterative process, each step of which consists in solving a simpler auxiliary extremal problem. Instead of the family of constraints, we use a new single constraint, which is a linear combination of initial constraints with the corresponding aggregated weight coefficients. After the auxiliary extremal problem is analytically solved, the aggregated coefficients are re-calculated again. In conclusion, we will illustrate some results of computational experiments.

## 3 Velocity reconstruction: proposed algorithm

Let us describe the proposed algorithm for solving the problem of reconstruction of velocity characteristics of a regional medium.

*Algorithm input data*: Domain $\Omega \subset \mathbf{R}^2$; coordinates of sources and receivers; travel times $t_i$, ($i = 1,...,n$); measurement errors $\varepsilon^{(1)}$, $\varepsilon^{(2)}$; initial velocity distribution $v^0(r)$, $r \in \Omega$; some partition of domain $\Omega$ into cells numbered $1,...,K$.

*Initial step of the algorithm*: From starting velocity $v^0(r)$, determine wave paths $l_{0i}$, and travel times $t_{0i}$, $i = 1,...,n$, i. e. solve the problem of ray tracing. For this purpose, apply one of methods of numerical integration of eikonal equations [21]:

$$\frac{dx}{ds} = vp_1, \quad \frac{dy}{ds} = vp_2, \quad \frac{dp_1}{ds} = -v^{-2}\frac{dv}{dx}, \quad \frac{dp_2}{ds} = -v^{-2}\frac{dv}{dy}.$$

Find time discrepancies $\delta t_i = t_i - t_{0i}$, $i = 1,...,n$. From rays $l_{0i}$, calculate matrix $A = \{a_{ij}\} \in \mathbf{R}^K \times \mathbf{R}^n$, whose rows $A_i \in \mathbf{R}^K$ are approximations of data kernels $G_i$, $i = 1,...,n$. While building the matrix, take into account the fact that theoretically kernels $G_i(x,y)$ must satisfy the relationships [1, 4, 11,18]:

$$\iint\limits_{\Omega} G_i(x, y)\,dx\,dy = t_{0i},$$

which may be transformed into the following:

$$\sum_{j=1}^{K} a_{jj}\Delta_S = t_{0i}, \quad (i = 1,...,n)$$

where $\Delta_S$ is the cell's square. In this case, the coefficients $a_{ij}$ are calculated by the formula:

$$a_{ij} = \begin{cases} \dfrac{|l_{0i}|}{v_{0i}n_i}, & \text{if } (x_j, y_j) \in O_h(l_{0i}), \\ 0, & \text{otherwise.} \end{cases}$$

Here, $|l_{0i}|$ is the length of ray $l_{0i}$; $v_{0j}$ is velocity at the $j$-th cell; $O_h(l_{0i})$ is an $h$-neighborhood of ray $l_{0i}$; $n_i$ is the number of cells, whose centers belong to $O_h(l_{0i})$.

Let us now describe an algorithm for solving a problem that generalizes the problem (4) in discrete formulation. The above calculations permit us to state the following optimization problem:

$$\begin{cases} \min (y'D' + c'y) \\ A^1 y \leq b \end{cases} \tag{5}$$

Here, $c,y \in \mathbf{R}^K$, $D \in \mathbf{R}^K \times \mathbf{R}^K$, $b \in \mathbf{R}^{2n}$, and prime denotes transposition. We can reduce the problem (5) to a particular case corresponding to the problem (4) assumed that

$$A^1 = \begin{pmatrix} A \\ -A \end{pmatrix}, \quad b = \begin{pmatrix} \delta t_i + \varepsilon^{(1)} \\ \delta t_i - \varepsilon^{(2)} \end{pmatrix}, \tag{6}$$

$$D = I, \quad c = 0,$$

where $I$ is the identity matrix.

After we substitute $z = Dy$, we have the following problem:

$$\begin{cases} \min \left| z^2 \right| + c'D^{-1}z \\ A^1D^{-1}z \le b. \end{cases} \qquad (7)$$

To solve the above extremal problem we apply the following modification of the algorithm suggested in [12]. We introduce parameter $\tau \in [0,1]$, $\tau_0 = 0$, number of steps $N$ of algorithm solving problem (7), and vector $x \in \mathbf{R}^{K+2n}$, consisting of two parts:

$$x = \begin{pmatrix} z \\ s \end{pmatrix}, \quad z \in \mathbf{R}^K, \quad s \in \mathbf{R}^{2n}.$$

Here we assume that $x_0 = 0$.

The *iterative process* includes the following calculations ($1 \le k \le N$):

$$x_{k+1}^j = x_k^j + u_k^j \frac{1}{N}, \quad j \in [1 : K + 2n],$$

$$\tau_{k+1} = (k+1)\frac{1}{N}.$$

Vectors $u_k \in \mathbf{R}^{K+2n}$ that play the role of control parameters are chosen according to the following rules:

$$u_k^j = \begin{cases} a_j, & d_{jk}^* < a_j \\ b_j, & d_{jk}^* > b_j \\ d_{jk}^*, & d_{jk}^* \in [a_j, b_j] \end{cases}, \quad j \in [1 : K],$$

$$d_k^* = -[P_k/\alpha + c^*/2],$$

$$P_k = E'(Ez_k + s_k - \tau_k b),$$

$$c^* = c'D^{-1}, \quad E = A^1D^{-1}, \quad j \in [1 : K];$$

$$u_k^j = \begin{cases} 0, & g_{j-K}^{(k)} \ge 0 \\ d_{j-K}, & g_{j-K}^{(k)} < 0 \end{cases}, \quad j \in [K+1 : K+2n],$$

$$g^{(k)} = E z_k + s_k - \tau_k b.$$

It is assumed that the value $\alpha$ is small enough.

*Algorithm output*: Vector $z_k = \{x_k^j\}_{j=1}^K$ is the solution of problem (7), which allows to find the solution of problem (5) by formula: $\hat{y} = D^{-1}z_k / \tau_k$. Under conditions (6), components of vector $\hat{y}$ are values of the discrepancy of slowness in cells: $\hat{m}_j = y_j$, $j \in [1:K]$.

**Remark.** Numerical simulations conducted so far indicate that in case of *N* big enough (near 100 and greater), it is sufficient to carry out only up to about *N*/3 steps of the algorithm. Further computations (*k* > *N*/3), as a rule, do not result in a significant increase of precision.

## 4. Experimental results



Fig. 2: The set of seismic rays under testing.

Efficiency of the algorithm was tested in the course of computational experiments. Let us consider an example of model seismic region under investigation. This example is quite indicative since the region contains inhomogenieties related to both higher and lower velocities. Fig. 2 shows the set of sources and receivers connected with 100 rays, which was used in our tests. The input data were travel times $t_i$, calculated for bent raypaths. In our experiments, the domain $\Omega$ was split into 40×40 square cells.

Initial (unknown) distribution $v_0(x,y)$ of signal propagation is shown in Fig. 3 as a function of spatial coordinates *x,y*. Background velocity in the region equals 3.0, velocities in the inhomogenieties being 3.7 and 2.5 respectively. Errors $\varepsilon^{(1)}$, $\varepsilon^{(2)}$ were assumed to be zero. Note that the implicit presence of errors in values $t_i$ is unavoidable due to inaccurate measurements or approximate calculations in case of artificial data.

Fig. 3: Model distribution $v_0(x,y)$ of seismic signal propagation velocity as a function of spatial coordinates.

Zero approximation of the velocity distribution was assumed to be uniform and equal to background velocity: $v^0(x,y) = 3$ for $(x,y) \in \Omega$; the initial value of the discrepancy was set to $m^0(x,y) = 0$ for $(x,y) \in \Omega$.



Fig. 4: Reconstructed velocity distribution.

The distribution of velocity reconstructed by the algorithm is shown in Fig. 4.

Consider also the distribution of value $\tilde{v}(r) = v(r) - v_0(r)$, i. e. difference of the reconstructed and the real velocities of the signal propagation. Let us call it the error of the velocity reconstruction. The distribution of $\tilde{v}(r)$, $r \in \Omega$, is shown in Fig. 5. As seen from the figures, the suggested algorithm gives rather sharp reconstruction of the field of velocity. Background velocity is identified quite precisely as well. It may be considered as an advantage that the algorithm gives not as smooth a solution as the previous algorithm described in [7].

The dependence between the average error of the velocity reconstruction over all cells and number of iterations is shown in Fig. 6. This figure confirms the remark given in the previous section concerning the reasonable number of iterations. As seen from Fig. 7, decrease of the number of iterations gives a significant gain in time.

Fig. 5: Distribution of error of velocity reconstruction produced by the algorithm.



Fig. 6: Dependence of average error of velocity reconstruction on number of iterations *N*.

80

Fig. 7: Dependence of computation time on number of iterations *N*.

## 5. Concluding remarks

In this paper we have introduced a new algorithm for the problem of reconstructing the property of the medium based on a set of experimental results characterized by measurement errors. While we have presented our data for the geoseismic problem a similar approach can be applied to the plasma physics problem described in Section 2. Both these problems are examples of application of computer diagnostics as tools for the development of methods and techniques in the risk-management area. In the case of nuclear physics, the correct assessment of various properties of the nuclear material is of extreme importance for the management of risk of nuclear plants. In the case of studies of structure of the Earth mantle, the produced data can be used in the context of risk assessment for the construction and insurance costs of various private and commercial structures (e. g. in California or other areas of high seismic activity). The proposed algorithm is shown to produce correct reconstruction of the property in question as well as being rather efficient.

As the next step we will apply the proposed approach to other similar problems, including real-world problems based on actual experimental data.

## References

1. Aki, K., Richards, P.G., 1980, Quantitative Seismology: Theory and Methods. Vol. 2. Freeman and Company, San Francisco.

2. Astafiev, N.N., 1991, Infinite systems of inequalities in mathematical programming. Nauka, Moscow (in Russian).

3. Blizorukova, M.S., Maksimov, V.I., 1997, On reconstruction of extremal input in a system with delay. *Vestnik PGTU. Functional-differential equations*, **4**, pp. 51–61 (in Russian) .

4. Backus, G., Gilbert, F., 1967, Numerical applications of a formalism for geophysical inverse problems. *Geophysical Journal of the Royal Astronomical Society*, **13**, pp.247–276.

5. Chernikov, S.N., 1984, Linear inequalities. Nauka, Moscow (in Russian).

6. Digas, B.V., Ermoliev, Yu.M., Kryazhimskii, A.V., 1998, Guaranteed optimization in insurance of catastrophic risks. IIASA Interim Report IR-98-082.

7. Digas, B.V., Maksimov, V.I., Bukchin, B.G., Lander, A.V., 1998, On an algorithm solving inverse problem of ray seismics, *Computational Seismology*, **30**, pp.207–224. (In Russian)

8. Ermoliev, Yu.M., Kryazhimskii, A.V., Ruszczyński, A., 1997, Constraint aggregation principle in convex optimization . *Mathematical Programming*, 76, pp. 353–372.

9. Ermoliev, Yu.M., Ermolieva, T., MacDonald, G., Norkin, V., 2000, Catastrophic risk management and economic growth, IIASA Interim Report IR-00-058.

10. Gabasov, P., Kirillova, F., Kostyukova, O., 1986, Constructive methods of optimization. Part 3: Net problems. Universitetskoe, Minsk (in Russian).

11. Jonson, E., Gilbert, F., 1972, Inversion and inference for teleseismic ray data. *Math. Comp. Phys.* V. 12. pp.231–266.

12. Kryazhimskii, A.V., 1999, Convex Optimization via Feedbacks. *SIAM J. Control Optimization,* vol. 37, pp. 278–302.

13. Kryazhimskii, A.V., Maksimov, V.I., 1998, On an iterative procedure for solving the control problem under phase constraints. *Computational Mathematics and Mathematical Physics*, Vol. 38, No. 9, pp.1484–1489 (in Russian) .

14. Kryazhimskii, A.V., Maksimov, V.I., Osipov, Yu.S., 1997,  On reconstruction of extremal disturbances in parabolic equations. *Computational Mathematics and Mathematical Physics*, Vol. 37, No. 3, pp.119–125 (in Russian)

15. Kryazhimskii, A.V., Maksimov, V.I., Samarskaya, E.A., 1997, On reconstruction of inputs in parabolic systems. *Mathematical Modeling*, **9** (3), pp.51–72 (in Russian).

16. Kryazhimskii, A.V., Osipov, Yu.S., 1987, On regularization of a convex extremal problem with inaccurately given constraints. An application to the problem of extremal control with phase constraints. *Some methods of positional and program control.* UB AS USSR, Sverdlovsk, pp.34–54 (in Russian).

17. Kuznetsov, E.I., Scheglov, D.A., 1980, Methods of diagnostics of high-temperature plasma. Atomizdat, Moscow.

18. Levshin, A.L., Yanovskaya, T.B., Lander, A.V., Bukchin, B.G., Barmin, M.P., Ratnikova, L.I., Its E.N. Surface seismic waves in horizontally non-uniform Earth. Nauka, Moscow, 1983. (In Russian)

19. Mourogov, V., Kagramanian, V., 2001, The case for innovative nuclear reactor and fuel cycle systems, Proceedings of the Workshop on Risk Management: Modeling and Computer Applications, IIASA, Laxenburg, Austria, 14–15 May, 2001;

20. Mourogov, V., Kupitz, J., 1998, Nuclear energy issues and the role of small and medium sized reactors, 23*rd* Annual Symposium of Uranium Institute, London, U.K., 9–11 September, 1998.

21. Psencik, I., 1994,  Seismic ray method for inhomogeneous isotropic and anisotropic media. Second Workshop on Three-dimensional Modeling of Seismic Waves Generation, Propagation and their Inversion. Italy, Trieste.

22. Tikhonov, A.N., Arsenin, V.Ya., Timonov, A.A., 1987,  Mathematical problems of computer tomography. Nauka, Moscow.

23. Vasiliev, F.P., 1980,  Numerical methods of solution of extreme problems. Nauka, Moscow (in Russian).

# Part III: Risk Management: Socio–Economic Aspects

# Innovation: risk and economic safety of the Ural region

*Marina Blizorukova, Andrei Maksimov, Andrei Shorikov*

## Abstract

The paper describes the innovation processes in the Urals, one of the most industrially oriented regions in the center of Russia. Some aspects of risk management in the region are discussed. Socio-economic aspects, as well as methods of risk assessment and the innovation activity in small business are considered.

# 1.  Introduction

The Russian economy interacts with the world economy closely enough nowadays. In this context, security of innovations and information protection in production in Russia become important not only for this country and its closest neighbors. Risky factors in innovation activities arise due to non-protection and misappropriation of technological information as well as by failures in decision-making and in the implementation of scientific and technological developments. The crisis of the Russian economy had led to the fact that the Russian industry has become not sensitive enough to the progress in knowledge-based technologies. The Federal Budget supports basic and applied research insufficiently, which leads to the fact that scientific achievements become not well secured by patents and laws. The system, which existed for years, has been destroyed. Enterprises have no prize funds to award inventors. Regulations on intellectual property are not well shaped, and the procedure of receiving patents (which includes preparing routine technical documents, manufacturing a sample, testing it, organizing an industrial process, etc.) is extremely complicated. All this has led to the deceleration of innovation processes in the regions of Russia, particularly, in the Urals. However, the necessity of stimulating the technological development in this region is obvious. Recently, the Regional Center for Control of Scientific, Technological and Innovation Activities was created in the Urals. The aim of the Center is protecting the authors' rights of patent holders.

Intellectual products, which have been kept for a long time on the interior Russian market only, have stable demand in countries with developed economies nowadays. This has induced the uncontrolled outflow of Russian technologies from Russia. In 1996-2000, more than a half of 500 applications of Russian researchers who have applied for patents in the USA have been registered, which contradicted the Patent Law of the Russian Federation. More than 90% of these applications have been sent to the Department for Patent and Commodity Marks of the USA. This example confirms that the indicated problem is important in the light of the international market of innovations.

Every innovation project assumes the participation of the following economical agents:

- The author of the project
- Investors financing the development and implementation of the project
- An enterprise
- Consumers of production (market)

Protection of the investment activity meets the interests of all agents listed above and promotes the successful implementation of the innovation project. It is also evident that the innovation activity makes sense only if the effectiveness of the expected result exceeds estimated risk. The higher is the effectiveness of the expected result, the higher is risk, up to which the agents may find it reasonable to start up.

The historical experience shows that the economies based on innovations win. The unstable situation in Russia is partially explained by the fact that the existing social institutes are unable to organize regular innovation activity in any concrete form manifested.

The non-competitive character of production in the soviet period and its defeat in the post-soviet time under onset of world importers on the Russian market were extreme economic perturbations; even the huge Russian structures specially oriented to promoting innovations were unable to hold the innovation potential under control. On the other hand, small industrial and commercial structures began to demonstrate competitiveness in innovations, although in local frames.

# 2.  Classification of innovations

Investments in material and nonmaterial assets connected with scientific and technological progress are usually called innovation investments. The innovation means introduction of a new

order, a new method, an invention, a new industrial sample, etc [3]. In a wider sense, innovations are understood as activities aimed at developing, creating and using new technologies, production, services, new organizational-technological and social-economical decisions on industrial, financial, commercial, or administrative character. The overall goal of innovations is always more complete satisfaction of human demands and the increase of profit on this base. The innovation process is a chain of events, through which the innovation is transformed from an idea to a concrete product, technology or service and extends into economic life [2]. Table 1 gives a classification of innovations with respect to different features.

**Table 1.**

| Feature | Type of innovation |
|---|---|
| 1. Degree of radicalism | Basic |
| | Improving |
| 2. Reason of rise | Reactive (a reaction for competitor's actions) |
| | Strategic (competitive advantage as a perspective) |
| 3. Contents and field of application | Product (new products, materials) |
| | Process |
| | Market innovation (new fields, new markets) |
| 4. Character of satisfying requirements | New requirements |
| | Existing requirements |
| 5. Scales of spreading | Creation of a new branch |
| | Application in all branches |
| 6. Degree of significance | Regeneration and adaptive changes of original features |
| | New variant |
| | New generation |
| | New form |
| | New sort |
| 7. Degree of novelty | On the base of a new scientific discovery (absolute novelty) |
| | On the base of a new method of application of phenomena discovered earlier (relative novelty) |
| 8. Stages of innovation process | Creation of innovation |
| | Development |
| | Distribution (replication) |
| | Usage |
| | Modification |
| | Retirement |
| 9. Stimulus of appearance (source) | Inspired by development of science and technique |
| | Inspired by requirements of industry |
| | Inspired by requirements of market |
| 10. Role in reproduction process | Consumer |
| | Investment |
| 11. Degree of complexity | Complex (synthetic) |
| | Simple |

## 3. Innovation cycle

Every innovation has a cycle, i.e., a time period, which starts with theoretical and applied research, includes stages of development and use of a new technological idea, improvements of its technical and economical parameters and other modifications, and terminates at a phase when this technique is substituted by a new and more effective one. The stages of the innovation cycle are shown in Fig. 1.



Fig. 1

Basic stages of an innovation cycle are, therefore, research and development (creation of innovations), mastering, serial (mass) production and servicing. Each stage can, in turn, be divided into sub-stages. For example, creation of innovations comprises developments in the areas of fundamental and applied research, as well as in the areas of design and experiment. Fundamental research is connected with extending, deepening and systematizing knowledge on a given subject resulting in the development of a new theory and new concepts. Experience shows that the investments in science yield, in the long run, a much higher profit than investments in other profitable spheres, for example, production.
Applied research having practical importance is connected with estimating possibilities of using the results of fundamental research to develop a particular branch, territory or industrial process. Machines, instruments, new technologies, etc. appear as results of applied research. Since a concrete result on this stage is not always evident, the investor must consider risk of losing part of his capital (investment risk).
The third stage of the innovation process is designing and experimental developments. This stage is aimed at manufacturing and testing samples of new products. On this stage, instructions and manuals for using new technologies and software are prepared.
The next stage is commercialization of the invented innovation and the demonstration of the innovation on the market. The main aim of the innovation activity is implementation of the innovation in industry. On this stage, one can find four temporal phases of a cycle: inculcation, increase of production and sales, deceleration of growth, and falling-off. On the stage of inculcation and demonstration, creation of industrial capacities, training the stuff, advertising

and other activities require large investments. On this stage, the reaction of market to the innovation is not clear yet.

The innovation process continues after the appearance of a new product (service) on market and reaching the project capacity in industry. On this stage, the innovation is improved and modified, and new customers, new markets and new fields of its utilization are looked for.


## 4. Risk and economic safety

Risk is the probability of higher losses in profit, compared to the expected ones. Making a decision on an investment, the investor should estimate risks of several kinds: risk of frustration of the project; risk of incomplete implementation; industrial risk; risk of failing in paying debts; taxation risk; market risk, etc. Ecological, criminal, political and other factors of risk should also be taken into account. Depending on the degree of risk, the investments can be classified in the following categories: investments with no risk, investments with acceptable levels of risk, investments with risk exceeding the acceptable level, and, finally, investments with critical or catastrophic levels of risk (Figure 2). Innovation risk never equals zero. It exists objectively, regardless of the people making the principal decision on starting off the innovation project, or people making local decisions during its execution. The people or organizations willing to adopt and implement the project must take into account the factor of risk. Not taking it into account may lead to over-normative reserves of unsold production, losses in profit, lowering the effectiveness below a planned level, ineffective expenses in material, labor and financial resources, losses in property, etc.



Fig.2

The scales of possible negative consequences of decisions made without taking into account risk factors can be large enough. The managers must take into consideration risks arising due to both making the principal decision on starting off the project and running the project. In the innovation activity, risk assessment is based on the scientific, technological, economic and sociological analysis of the project as a source of risk, and also on the analysis of external and internal risk factors, estimation of acceptable levels of risk and creation of mechanisms and

models interconnecting risk indices and factors. The innovation aims at increasing the economic potential of an enterprise. The chosen index of risk level must, therefore, determine admissible ranges of deflection in the increase of the enterprise's potential and also admissible damages due to these deflections. One can define the following factors of innovation risk:

- wrong interpretation of the aims determining the development of the innovation project;
- errors in the estimation of the economic potential (the origin of this risky factor may be stipulated by inaccuracies in data and the lack of reliable information on future technological changes;
- loss or misrepresentation of information during the transition to tactical planning.

Conditionally, the risk factors can be classified into external and internal ones [1]. External risk. factors are not connected with the activity of the enterprise realizing the innovation project; usually, these factors lie in the areas of politics and country's economy. External risk factors lying beyond these areas are also highly important for the modern innovation activity. In Russia, they are stipulated by regional conditions and include ecological risks (pollutions), social risks (refugees, unemployment, low incomes of population), local economic risks (non-self-sufficiency of the territory in basic products, limited labor resources, non-satisfactory economic-geographical position, poorly developed infrastructures, etc.)

The internal risk factors are industrial risk factors such as violation of technological discipline by the staff, unscheduled breaks in the technological process; damages; violation in delivering raw materials and details, economic crime, and others. Taking into account external and internal risk factors leads to two directions in using the concept of acceptable risk. First, it is used to order innovation projects and select those acceptable for a concrete enterprise. Second, it is used for the development of measures on lowering risk in the process of the implementation of the selected project. In planning such measures, the issue of economic security of the innovation activity plays an important role. Final strategic decisions in the innovation activity of an enterprise should be strictly controlled so as the long-term goals of the enterprise connected with the improvement of its financial state and growth in scales of its management activity will not be overshadowed by intermediate needs of different departments.

Since any innovation process is a result of a set of decisions, it is important to work out a sequence of risk estimates given by different participants of the innovation process. The lack of coordination in estimating risk by the author, investor and executor leads, as a rule, to the frustration of the project. Every innovation project needs to be clearly described. To achieve this, one can recommend using partially structured description forms (business-plans) based on the well-known form of business-plans for enterprises. However, this scheme must have a number of specific features directed to the more precise estimation of risk. These features essentially depend on the type of the project.


## 5. Methods of analysis of risk

Different methods can be used for the quantitative analysis of risk. Nowadays, widely used are statistical analysis, analysis of advisable expenditures, the method of expert estimations, analytic methods, usage of analogues, etc. We dwell upon the first two methods.

Statistical methods of risk analysis consist in the following. For the calculation of the probabilities of losses, all statistical data connected with the effectiveness of financing innovation projects are investigated. The frequency of arising of some level of losses indicates the probability of entering one of the following risk domains: domain of no risk, the domain of acceptable risk, domain of risk higher than acceptable, and domain of critical (catastrophic) risk. The frequencies are given by the formula

$$f_i^0 = \frac{n_i}{n_{tot}}$$

where $n_i$ is the number of occasions of entering risk domain $i$ according to the statistics, and $n_{tot}$- the total number of occasions in the statistical selection. Using such calculations, one estimates about the degree of risk during the implementation of the project.

The Monte-Carlo technique is a popular method of statistical experiments. An advantage of this method is that it gives a possibility to analyze different risk factors in the framework of one approach. Different types of projects differ in their vulnerability with respect to risks.  A disadvantage of the statistical risk analysis is the use of probabilistic characteristics, which makes it difficult to apply the results of the analysis in practical activities. The analysis of advisable expenditures is oriented to identification of possible risk zones. Overexpenditures can be stipulated by one of the four principal factors or their combinations: underestimation of initial costs; change of limits of design; differences in productivities; increase of the initial cost. These principal factors can be worked out in details. On the base of a typical enumeration one can make a detailed control enumeration for a concrete project or its elements. There is a possibility to minimize capital exposed to risk by means of dividing the process of investment into stages. On each stage, the investor taking into account the analysis current risk is allowed to terminate investments.

The estimation of the degree of risk via the method of advisable expenditures uses three indices of financial stability of the project: surplus or deficiency in internal funds ($E^C$); surplus or deficiency in own, medium- and long-term loan sources of forming reserves and expenses ($E^T$); surplus or deficiency in total value of basic sources of forming reserves and expenses ($E^H$). These indices correspond to the indices of supply of reserves and expenses.

To identify a financial situation, one can use the following three-component index:
$$S=\{S(E^C)S(E^T)S(E^H)\}$$
where
$$S(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$$
Depending on the value of $S(x)$ the project under consideration may be attributed to the following risk domains:

1.  Domain of no risk:
$S=(1,1,1),\quad E^C{\geq}0,\quad E^T{\geq}0,\quad E^H{\geq}0$

2.  Domain of acceptable risk:
$S=(1,1,1),\quad E^C{=}0,\quad E^T{=}0,\quad E^H{=}0$

3.  Domain of risk higher than acceptable:
$S=(0,1,1),\quad E^C{<}0,\quad E^T{\geq}0,\quad E^H{\geq}0$

4.  Domain of critical risk:
$S=(0,0,1),\quad E^C{<}0,\quad E^T{<}0,\quad E^H{\geq}0$

5. Domain of catastrophic risk:
$S=(0,0,0),\quad E^C{<}0,\quad E^T{<}0,\quad E^H{<}0$

## 6. Innovation activity and small business in the Ural region

In recent years the small business is more and more actively attracted to the innovation activity in the Sverdlovsk region, the Urals. The main activities of small enterprises, the structure of their economic characteristics, and the structure of employment in small business are presented in the table and figures below.

Among the small business enterprises actively taking part in the innovation activity there is a number of «business-incubators» and «scientific and technological parks» of the Sverdlovsk region. For example, the business-incubator of «The Ural technological park» stretching for 1200 square meters includes more than 20 innovation firms. One of the most important organizations is the non-commercial company «Innovation Technological Center »Academical» found by the government of the Sverdlovsk region, the Ural Branch of the Russian Academy of Sciences and the Federal Fund for Supporting Development of Small Enterprises in the Scientific-Technological Sphere. The main tasks of the Center are creation and development of the infrastructure of small innovation business, implementation of new technologies and help in creation of new vacancies for the effective use of scientific-technological potential. On the territory controlled by the Center there are small enterprises producing anticorrosion covering, silicon and organic combinations, and utilizing waste products to extract silver and non-ferrous and precious metals [4].

In the framework of the technological park «Visokogorskii» nine projects for environment's protection and ecology have been started. Special attention is given to the conversion of waste products of metallurgy, land recultivation, and providing resource and energy savings. The projects are cunducted within the framework of federal and regional programs on energy and ecology and supported by enterprises, the Federal Budget, and Regional Ecological Fund. Scientific research of applied character is carried out by regional scientific institutions, including the Ural State Geological Academy, the Ural Institute of Metals and others. The regional scientific-technological park «Uralskii» is created in the Ural State Technical University. It comprises 38 enterprises of small science-consumable business. The principal directions of their activity are development and implementation of new technologies in radio-engineering, electronic communication engineering, instrument-making, chemical industry and medical technique.

**Structure of employment in small enterprising in the Sverdlovsk region in 2000** [8]

**Employment by types of small enterprising**



**Employment in small enterprises: classification in sectors**

Percentage of permanent employees:

| | |
|---|---|
| Commerce and public catering | 32,9 |
| Industry | 22,8 |
| Building | 19,2 |
| Science | 3,4 |
| General business activity | 2,9 |
| Transport | 2,7 |
| Public health and sports | 2,6 |
| Culture and education | 1,0 |
| Other | 12,5 |

**Employment in small enterprises: classification in kinds of property**

- private enterprises                                      93,5 %
- enterprises of the combined form of property
  (no foreign participation)                            5,5 %
- enterprises of the combined form of property
  (Russian and foreign participation)              1,0 %

**The number of workers in small enterprises** [8]



**others**
**22%**

**science and**
**scientific service**
**3%**

**building**
**19%**

**industry**
**23%**

**commerce and**
**public catering**
**33%**

**The number of workers in small enterprises: 140900.**

**Branch structure of economic indices of small enterprises of the Sverdlovsk region in 2000 [7]**

## 1. Volume of produced goods and services

others
21%

science and scientific service
6%

industry
30%

building
19%

commerce and public catering
24%

## 2. Profit

others
25%

industry
20%

science and scientific service
17%

building
14%

commerce and public catering
24%

**Economic indices of development of small enterprises of the Sverdlovsk region**

**Main indices of development of small enterprising of the Sverdlovsk region in 1998–2000 [6]**

| INDICES | 1998 | 1999 | 2000 |
|---|---|---|---|
| 1.  Amount of types of small enterprising: | | | |
| • Small interprises (units) | 31381 | 25900 | 31891 |
| • Individual enterprisers (people) | 110900 | 109550 | 92875 |
| • Farms (units) | 2359 | 2162 | 2121 |
| 2.  Amount of small enterprises in branches of economy: | | | |
| • Industry and building (units) | 9854 | 8500 | 10163 |
| • Commerce and public catering(units) | 15102 | 11450 | 14398 |
| • Science and scientific service (units) | 1232 | 1400 | 1549 |
| • General business activity (units) | 1161 | 1000 | 1183 |
| • Others (units) | 4032 | 3550 | 4598 |
| 3.  Permanent employees in the sphere of small enterprising (thousands people): | 270,3 | 263,8 | 241,4 |
| • At small enterprises (thousands people) | 151,9 | 146,9 | 140,9 |
| • Individual enterprisers (thousands people) | 110,9 | 109,5 | 92,9 |
| • Farmers (thousands people) | 7,5 | 7,4 | 7,6 |
| 4. Part-time job (thousands people) | 34,9 | 29,5 | 29,3 |
| 5. Percentage of small enterprising in region economy: | | | |
| • volume of industrial production | 2,3 | 2,8 | 2,8 |
| • volume of building-assembly works | | | |
| • retailing commodity circulation | 24,1 | 23,8 | 24,2 |
| • volume of farm's products | 22,4 | 27,5 | 25,5 |
| • permanent employees | 3,2 | 1,0 | 1,2 |
| | 13,2 | 13,5 | 12,4 |

## References

1. Control of investments and innovations, 1997.

2. Hitman, L.D., Jonk, M.D., 1997, The basis of investment.

3. Dones, D., Goodman, D., 1997, Financial-investment vocabulary.

4. Web site «Science in regions»
   (**http://www.extech.ru/s_e/region/subjects/sverdlov/index.htm**)

5. Web site «Export map of the Sverdlovsk region»
   (**http://www.midural.ru**)

6. Main indices of development of small enterprisies of the Sverdlovsk region from 1998 till 2000 years». Sverdlovsk region committee of state statistic.

7. Information of the Sverdlovsk region committee of state statistic and municipal authorities.

8. Express information «Main indices of development of small enterprisies of the Sverdlovsk region», 2000. Sverdlovsk region committee of state statistic.

# Financial risk management: set–valued uncertainty modeling

*Oleg Nikonov*

## Abstract

Traditional approaches and methods used in control problems motivated by financial management are based on the probability theory and stochastic calculus. Well known fundamental results have been obtained in this field during the last three decades. But for dealing with emerging markets, with their lack of statistics and high level of uncertainty, the approaches based on set-valued uncertainty models seem to be adequate and useful tools. In the paper we consider a problem of dynamic investment portfolio selection treated via guaranteed control theory. A formalized setting and solution that combine the methods of this theory with traditional mean-variance approach are given.

# 1. The Markowitz–Tobin static portfolio selection problem

A classical version of the portfolio selection problem presumes the presence of N risky assets (securities) with specified returns (rates of return) $r_i$ ($i = 1, \ldots, N$), which are treated as stochastic variables. Denote by $x_i$ and $V = \{\sigma_{ij}\}$ their mean values and covariance matrix respectively.

It is also assumed that there is an opportunity of risk-free investment. Let $x_0$ denote the corresponding rate of return. The portfolio is associated with an $(N+1)$-vector

$$\hat{y} = (y_0, y)^T = (y_0, y_1, \ldots, y_N)^T \in R^{N+1}, \quad \sum_{i=0}^{N} y_i = 1,$$

each coordinate of which corresponds to a share of capital invested into the related financial instrument (symbol ''$T$" stands for transposing).

Values $y_i$ are not supposed to be nonnegative, which reflects the opportunity of borrowing, lending, or short selling.

The expected portfolio's rate of return is then defined by relation $\mu(\hat{x}, \hat{y}) = y_0 x_0 + (y, x)$, where $\hat{x} = (x_0, x)^T$, and standard deviation $\sigma(y) = (y^T V y)^{1/2}$ is interpreted as the corresponding risk. Thus, with each portfolio $\hat{y}$ a point on the plane $\sigma, \mu$ (risk-return) can be associated. A portfolio $\hat{y}^* = (y^*_0, y^*)^T$ is called efficient or nondominated if there is no portfolio $\hat{y}$, for which $\mu(\hat{y}, \hat{x}) \geq \mu(\hat{y}^*, \hat{x})$ and $\sigma(\hat{y}) \leq \sigma(y^*)$, where at least one of the inequalities is strict.

The problem of constructing the set of efficient portfolios consisting of the pure risk assets has been posed and solved by H. Markowitz (1952). The solution is well-known and (under conditions that the matrix $V$ is positive definite and the vector $x$ has at least two different components) can be easily obtained by the Lagrange multiplier method. We will further assume the above conditions to be satisfied. Geometrically, the nondominated portfolios can be figured on the plane $\sigma, \mu$ by a branch of a hyperbola, the equation of which is derived in explicit form. The presence of a risk-free asset just simplifies the solution (Tobin, 1958). Here the relation connecting the portfolio's risk and its rates of return is linear: $\mu = x_0 + g\,\sigma,$

$$g^2 = (x - x_0 e)^T V^{-1} (x - x_0 e), \qquad e = (1, \ldots, 1)^T \in R^N.$$

The vector of the risk part of an efficient portfolio, which corresponds to the given rate of return $\mu$ ($\mu \geq r_0$), or to the given level of risk $\sigma$, can be expressed in the following way:

$$y = V^{-1}(x - x_0 e) g^{-2} (\mu - r_0)$$

or

$$y = V^{-1}(x - x_0 e) g^{-1} \sigma,$$

respectively.

The last relations reflect, in particular, the fact that the structure of the risk part is the same for all efficient portfolios. This fact is a basic one for the Capital asset pricing model (CAPM): if all investors follow the risk-return criterion in the above sense, each of them should form his or her portfolio distributing the risky assets in the same proportions. The optimal solution in the framework of this model is the distribution of funds between two portfolios-risk-free and fixed structure risky. The mentioned result is often called the two-fund theorem. From this, in turn, the property is deduced that the named structure must coincide with that of the real market portfolio, and the equilibrium prices and various characteristics of securities are determined.

## 2. Dynamical portfolio reconstruction: methods of the guaranteed control theory

Relations that determine the evolution of parameters $x$ and $V$, which characterize stochastic asset returns can be represented by stochastic differential equations. This approach has been realized, in particular, in (Merton, 1973), where the equations for the vector $x$ were taken in the form

$$dx_i = a_i dt + b_i d\eta_i, \qquad i=1, \dots, N,$$

and $\eta_i$ are independent standard Wiener processes.

We consider the following problem statement, assuming also that $x_i$ vary with time, but their dynamics is defined by the differential inclusion

$$\frac{d\hat{x}}{dt} \in A(t)\hat{x} + Q(t), \quad t_0 \leq t \leq \theta, \tag{1}$$

$$\hat{x}(t_0) = \hat{x}^0. \tag{2}$$

The multivalued function $Q(t)$ reflects uncertainty in the variation of expected returns. We shall assume the mapping $Q(t)$ to be piecewise continuous in t, the sets $Q(t)$ to be convex and compact, and $0 \in Q(t)$. Matrix $V$ is still supposed to be constant and positive definite.

Suppose that one can change the portfolio structure at each instant of time $t \in [t_0, \theta]$ with bounded velocity. In this case the portfolio dynamics is described by an equation $dy/dt = u$ with $y_0 = 1 - (e, y)$, where $u$ is a control function restricted by an inclusion $u \in P(t)$; $P(t)$ is assumed to have properties similar to those of $Q(t)$.

The aim of control is to ensure portfolio efficiency and to provide its specified characteristics according to the criteria of return and risk (mean - variance). The problem statement presumes construction of a control strategy realized in the form of control action $u=u(t)$, which is formed on the basis of information available at time t, and guarantees the required properties of the portfolio. The strategy is a rule that determines the dynamics of capital shares invested in risky and risk-free assets.

Formally, we define an admissible control strategy as a multivalued mapping $U=U(t, \hat{x}, \hat{y})$, which is measurable in t and upper semi-continuous in $\hat{x}, \hat{y}$ with convex compact values $U(t, \hat{x}, \hat{y}) \subseteq P(t)$.

Dynamic portfolio reconstruction is formalized as follows. Inclusions (1) and

$$\frac{dy}{dt} \in U(t, \hat{x}, \hat{y}) \tag{3}$$

under the above assumptions have an absolutely continuous solution $\hat{x}(t)$, $\hat{y}(t)$, where $y_0(t) = 1 - \square(e, y(t))$, $t_0 \leq t \leq \theta$ for any initial conditions (2) and

$$\hat{y}(t_0) = \hat{y}^0, \quad y_0^0 = 1 - (e, y^0).$$

Any solution can be prolonged on the whole interval $[t_0, \theta]$.

For each solution $\hat{x}(t) = (x_0(t), x(t))^{\mathrm{T}}$, $\hat{y}(t) = (y_0(t), y(t))^{\mathrm{T}}$, one can consider the evolution of portfolio risk-return characteristics:

$$\mu(t) = y_0(t)x_0(t) + (y(t), x(t)),$$

$$\sigma(t) = (y^T(t)V^{-1}y(t))^{1/2}$$

On the plane $\sigma$, $\mu$, a moving straight line corresponds to the set of efficient portfolios. To ensure the portfolio efficiency, one has to provide the following relation between the values $\mu$ $(t)$ and $\sigma(t)$:

$\mu$ $(t) = x_0(t) + g(t)$ $\sigma$ $(t)$,

where

$$g(t, \hat{x}) = \sqrt{(x(t) - x_0(t)e)^T V^{-1}(x(t) - x_0(t)e)} = \left\| x(t) - x_0(t)e \right\| V^{-1}.$$

The aim of control is, in particular, to guarantee the efficiency of the portfolio throughout the whole time interval $[t_0, \theta]$ provided that it has this property at the initial instant $t = t_0$.

We consider the case when the risk-free rate of return is constant $x_0(t) = r_0$, $dx_0/dt = 0$, and to simplify the formulae , we assume that $A(t) \equiv 0$. This does not lead to the loss of generality, at least theoretically, due to existence of an appropriate coordinate transformation. Hence,

$$\mu(\hat{x}(t), \hat{y}(t)) = y_0(t)r_0 + (x(t), y(t)), \quad \sigma(y(t) = \sqrt{(y^T(t)Vy(t)} = \left\| y(t) \right\| V.$$

Dynamic portfolio efficiency is guaranteed by the equality $\mu(\hat{x}(t), \hat{y}(t)) = r_0 + g(t)$ $\sigma(y(t))$.

Let $\hat{y}^0$ be an efficient portfolio for a given $\hat{x}^0$ with expected return $\mu^0 = \mu^0(\hat{x}^0, \hat{y}^0)$ and risk $\sigma^0 = \sigma(y^0)$.

Consider the following problems.

*Problem 1.* Specify an admissible control strategy $U = U(t, \hat{x}, \hat{y})$ that guarantees the efficiency of the portfolio $\hat{y}(t)$ for $\hat{x}(t)$ $(t_0 \le t \le \theta)$, whichever solutions $\hat{x}(t)$, $\hat{y}(t)$ to the differential inclusions (1)-(4) are taken.

*Problem 2.* Specify an admissible control strategy that solves Problem 1 and ensures a prescribed level of risk $\sigma(y(t)) \le \sigma^*$ and a prescribed level of expected return $\mu(\hat{x}(t), \hat{y}(t)) \ge \mu^*$.

Here $0 < \sigma^0 \le \sigma^*$ and $\mu^* \ge \mu^0 > r_0$.

The following regularity condition is assumed to be satisfied:

$$\left\| x^0 - er_0 \right\| V^{-1} - \max_{\|l\|V^{-1}} \int_{t_0}^{\theta} \rho(l \,|\, V^{-1}Q(t))dt = d > 0,$$

where $\rho(l \,|\, Q(t))$ is the support function of the set $Q(t)$, $\rho(l \,|Z) = max\{(l,z) \,|\, z \in Z\}$.

*Lemma 1.* The inequality (5) ensures the relation

$$\left\| x(t) - r_0 e \right\| V^{-1} \ge d$$

for all $t \in [t_0, \theta]$, which means the impossibility of situation when all rates of return of risky assets are simultaneously close to the risk-free interest rate.

Solvability conditions for Problems 1-2 are determined by the relations between the multivalued mappings $Q(t)$ and $P(t)$. It can be shown that a sufficient condition for the solvability of Problem 1 is the inequality

$$d \cdot \rho(l \mid P(t)) - \|y(t)\| V^{-1} \rho(l \mid V^{-1} Q(t)) \geq 0, \tag{7}$$

which should hold along possible solutions $y(t)$ of (1)-(4).

*Lemma 2.* If for an admissible control strategy $U = U(t, \hat{x}, \hat{y})$ the relation (7) holds for each solution to (1)-(4), then the strategy solves Problem 1.

Conditions formulated in terms of parameters initially given can be derived for both Problem 1 and Problem 2 on the basis of estimates for solutions $y(t)$ and $x(t)$.

Denote $\gamma = \max_{\|q\| V = 1} \rho(q \mid Q(t))$ and consider the following relations.

$$\kappa \cdot \rho(l \mid P(t)) - \sigma^*[\rho(l \mid V^{-1} Q(t)) + \gamma] \geq 0,$$

$$\forall l : \|l\|_V^{-1} = 1, \quad \forall t : t \in [t_0; \theta], \tag{8}$$

$$d^2 \cdot \rho(l \mid P(t)) - (\mu^* - r_0)[\rho(l \mid V^{-1} Q(t)) + \gamma] \geq 0,$$

$$\forall l : \|l\|_V^{-1} = 1, \quad \forall t : t \in [t_0; \theta], \tag{9}$$

where, $\kappa = \dfrac{(\mu^* - r_0)}{\sigma^*}$

$$\kappa \leq d. \tag{10}$$

*Theorem 1.* Under conditions (5) and (8)-(10) the set of efficient portfolios $\{\hat{y}\}$, which satisfy the relations $\mu(\hat{y}, \hat{x}(t)) \geq \mu^*$ and $\sigma(y) \leq \sigma^*$ is nonempty for each $t \in [t_0, \theta]$, and there exists an admissible strategy which solves Problem 2.


## 3. Numerical results

In conclusion, we describe some results of numerical simulation of the control procedures based on the methods presented in the paper. As input parameters, real data on Russian and foreign financial markets were used. In figures below, the results of modeling are represented for the six Russian blue chips: RAO EES, LukOil, Sberbank, Rostelekom, Surgutneftegas, Mosenergo. The period of time from October 1998 to August 1999 has been taken.

Fig. 1 reflects the dynamics of return for each security under consideration. The straight line is the return of the portfolio constructed according to the solution for Problem 2, where it was taken $\mu^* = \mu^0$, $\sigma^*$ is supposed to be sufficiently large, and which is therefore intended for guaranteeing a prescribed level of return without any restrictions on risk.

In Fig.2, the corresponding risk evolution is presented, which in this case changes considerably with time. Fig.3 depicts the risk evolution, when the problem of ensuring its given level (Problem 2 with $\sigma^* = \sigma^0$ and without any restrictions on return ) is solved.

Lastly, in Fig.4, the graph of the so-called stability index is presented. It is nothing else than a quantitative estimate of control resources necessary to solve Problem 2 in the previous particular case. The term ''control resources" we treat here as a size of the set $P(t)$, which has been taken as a ball with a constant radius (the straight line in the picture). This a priori chosen radius ensures the inequality (9) in which the set $Q(t)$ is replaced by the really observed values

of uncertain parameters. The graph of the corresponding value evolution is situated below the above mentioned straight line, which reflects the fact that the solvability conditions for the problem are fulfilled.



Fig.1. Securities and portfolio returns

Fig.2. Portfolio risk (1)



Fig.3. Portfolio risk (2)



Fig.4. Stability index

## 4. Conclusion

In the paper we consider a problem of dynamic reconstruction of an investment portfolio. The aim of control is to provide the portfolio efficiency within a given time interval in the mean—variance sense. A formalized setting and solution based on the approaches of the guaranteed control theory and game-theoretical methods are given. Theoretical results are illustrated with computer simulations involving real data of the Russian financial market. The results obtained seem to be useful while modeling emerging markets, characterized by the lack of statistics and high level of uncertainty.

## References

1. Markowitz, H., 1952. Portfolio Selection. *J. Finance,* 7, pp.77 – 91.

2. Tobin, J., 1958, Liquidity Preference as Behavior Towards Risk., *Rev. Economic Stud.*

3. Krasovskii, N.N. and Subbotin, A.I., 1988, Game-theoretical control problems. Springer, Berlin.

4. Krasovskii, A.N. and Krasovskii, N.N., 1995, Control under Lack of Information. Boston.

5. Kurzhanskii, A.B., 1977, Control and observation under uncertainty conditions. Nauka, Moscow. (in Russian)

6. Kurzhanskii, A.B., and V\'alyi, I., 1996, Ellipsoidal Calculus for Estimation and Control. Boston: Birkh\"auser.

7. Kurzhanskii, A.B. and Nikonov, O.I., 1990, On the Problem of Synthesizing Control Strategies: Evolution Equations and Multivalued Integration. *Sov. Math. Dokl.,* 41, N 2, pp.300 – 305.

8. Kurzhanskii, A.B. and O.I. Nikonov, O.I., 1994, Evolution Equations for Bundles of Trajectories of Synthesizing Control Systems, *Russ. Acad. Sci. Dokl. Math.,* 48, N 3, pp.606 – 611.

9. Samuelson, P.A., 1965, Rational Theory of Warrant Pricing, Industr. Manag., 13 – 31.

10. Merton, R.C., 1973, An intertemporal capital asset pricing model. Econometrica, 41, N 5, pp.867 – 886.

11. Bachelier, L., 1900, Theorie de la Speculation. Ann. de l'Ecole norm. super., 3.

12. Nikonov, O.I., 1998, Financial Decisions via Methods of Guaranteed Control Theory. Pliska. Stud. Math. Bulgar., 12, pp.133 – 140.

# Part IV: Risk Management: Optimizational and Environmental Aspects

# The great Caspian gas pipeline game

*Ger Klaassen, R. Alexander Roehrl[1], Alexander Tarasyev[2]*

## Abstract

In this paper we consider the problem of competition between gas pipeline projects. This problem becomes especially important in the context of evaluation of the future energy infrastructure, transportation routes, investments, timing of the projects, supply and consumption, price formation, etc. We illustrate significance of these aspects in the case study of the planned Caspian gas pipeline routes to the Turkey energy market. We propose a dynamical game model for description of investment scenarios, optimization of commercialization times of pipelines, regulation of gas supply and formation of gas price. This model constructed on the basis of classical micro and macro patterns of mathematical economics provides a macroeconomic tool for the analysis of future gas infrastructures. It comprises four microeconomic levels of optimization: assessment of the market of potential innovations, selection of innovation scenarios, regulation of the future supply and optimization of the current investments. The first simulations of the model give promising results for assessment of energy markets.

# 1. Introduction

The routing of oil and gas pipelines in Asia and especially the Caspian regions is at the center of the geopolitics of energy. Various countries in the Caspian region are playing the pipeline game to get access to one of the most promising markets in the region: Turkey. Turkey's gas demand is expected to quintuple by 2010 (EIA, 2000). Russia's Gazprom proposes to build the "Blue Stream" pipeline under the Black Sea to expand its current gas deliveries to Turkey. Technical difficulties might force Gazprom to look for alternative routes via Bulgaria or Armenia.. Turkmenistan, backed by the USA, is heading for the Trans-Caspian gas pipeline to deliver gas to Turkey. This pipeline would flow underneath the Caspian Sea through Azerbaijan and Georgia on to Turkey (EIA, 2000). Meanwhile the Iranians have completed their own gas pipeline to the Turkish border and are awaiting the Turkish side of the pipeline to be completed (see [Ignatius, 2000]). It seems that some of countries are moving ahead fast so as to preempt the investments decisions of others making it unattractive to built a new transmission pipeline since the market is not big enough. In addition gas could and actually is being shipped to Turkey in the form of LNG from Algeria and Egypt.

To increase the complication, Turkey is not the only relevant gas market for the three suppliers. Gazprom has expressed a desire for diversification of export routes. Gas market liberalization in Europe is expected to lower prices which lowers the rate of return of large scale investments in the Yamal peninsula or the Barents Sea. Gazprom is actively looking for new markets in the Asia Pacific (China) and the Middle-East (Turkey) (see [Makarov, 1999]) and several pipelines have been proposed to China. Turkmenistan could also deliver gas to China or to India and Pakistan (see [Klaassen et al., 2001]). Iran is looking into the option to pipe gas to India as well.

In order to analyze the indicated problems of development of the energy infrastructure we propose a new mathematical model of competition of large-scale gas pipeline projects. This model is constructed on the basis of classical micro and macro patterns of mathematical economics (see [Arrow, Kurz, 1970], [Intriligator, 1971]) and combines microeconomic levels of optimization: assessment of the market potential innovation, selection of innovation scenarios, regulation of the future supply and optimization of the current investments, through the price formation mechanism into macroeconomic tool for decision making. At all four levels of the model, dynamic optimization principles of optimal control and differential games (see [Pontryagin et al., 1962], [Krasovskii, Subbotin, 1988]) in construction of feedback solutions are significantly employed. The model takes into account the stages of construction and exploitation of the pipelines. At the stage of exploitation, as gas supply policies compete on market, decision making is relatively clear: the competitors search for an equilibrium supply at any instant. At the  stage of construction investment policies compete in the area of commercialization times and decision making is concerned with strong long-term aftereffects. A crucial decision is the choice of the commercialization time of the project, i.e., the time of finalizing the construction of the pipeline. Elements of optimal timing models (see [Barzel, 1968], [Tarasyev, Watanabe, 2001]) constitute the basis for a reasonable optimization setting in terms of commercialization (innovation) times. In the model we consider the interaction of competitors controlling the process by selection of commercialization times, current investments into pipelines and future supply within the framework of the game theory, (see, for example, [Basar, Olsder, 1982]). A right individual choice of the commercialization time could be the best response to the choices of the other competitors or a reasonable analogue of the best reply dynamics (see [Hofbauer, Sigmund, 1988], [Kryazhimskii et al., 2001]). Accordingly, a collection of commercialization times is suitable to every competitor if and only if the commercialization time of every competitor responds best to the commercialization times of the other competitors. Such situations can be associated with Nash equilibria in the game between the competing gas pipeline projects.

## 2. Gas trade models in the literature

| Type of Model | Macroeconomic Gas Trade Models | | |
|---|---|---|---|
| **Name** | Golombek et al. (1995) | GASTALE | GTM |
| **Organization** | Frischsenteret, Oslo | ECN | Stanford |
| **Source** | Golombek et al., (1995) | Van Oostvorn and Boots (1999) | Beltramo, Manne, Weyant, (1986) |
| **Geographical Coverage** | Western Europe (12 countries) | 7 largest gas consumer countries in EU | Canada, US, Mexico; 11 supply and 14 demand regions |
| **Time Scale** | 1990-2010 | 1995-2010 (to be extended to 2020) | 1990-2000 |
| **Software, Methodology** | GAMS, gas and electricity markets; profit maximizing Cournot producers facing an ideal third party access regime for gas transport; final demand for gas, oil, coal and electricity from two sectors ("households" and "firms"); Endogenous:all energy prices, used and traded quantities, and production of electricity. Exogenous: supply of natural gas. | Motivation: modeling EU gas directive (full competition and semi-open scenario); two liberalization steps (2000-2005 and 2005-2010); production companies and consumers but not distribution modeled explicitly; oligopoly -> Cournot equilibrium; three sectors: households, industry and power; drawing heavily on Golombek's model. | GAMS, market equilibrium model; US market deregulated and Canada and Mexico maintain export controls; GNP and price of oil exogenous; (max(profits)). |

Table 2-1 Overview of the Literature: Macroeconomic Gas Trade Models

| Type of Model | Game-Theoretic Models | |
|---|---|---|
| **Name** | DYNOPOLY | Wolf and Smeers, (1997) |
| **Organization** | Statistics Norway | Universite de Louvain et de Lille |
| **Source** | Brekke et al., (1987, 1991). Berg et al., | Wolf and Smeers, (1997) |

| | (1998). | |
|---|---|---|
| **Geographical Coverage** | external supply from Russia, Algeria and Norway; demand for 13 WEU countries | Norway, CEI, NL, Algeria, UK |
| **Time Scale** | 1995-2075 | 1990-2000 |
| **Software, Methodology** | Dynamic oligopoly (Cournot) model of gas supply; connected through gas price to exogenous demand input from SEEM model; Bertrand game for three producers; gas price relative to oil and coal price which are exogenous; Motivation: impacts of European $CO_2$ emissions of a reduction in Norwegian gas sales. | Stochastic version of a Stackelberg-Nash-Cournot Equilibrium Model; supply-side: oligopolistic market, followers in Nash equilibrium (max(profits)); demand side: uncertainty for the leader. |

Table 2-2 Overview of the Literature: Game-Theoretic Gas Trade Models

| **Type of Model** | **Cost-Optimization Models** | |
|---|---|---|
| **Name** | CPE | Energy Infrastructure Model for Asia/Eurasia |
| **Organization** | CRIEPI | University of Tokyo |
| **Source** | Sugiyama (1999) | Fujii and Yamaji (1999) |
| **Geographical Coverage** | 29 Chinese provinces | 37 energy "nodes" in China, Korea, Taiwan and Japan; plan to "extend geographical coverage" |
| **Time Scale** | 1990-2030 | 2000-2050 |
| **Software, Methodology** | Motivation: quantitative statements about acidification in China (SOx, Nox and $CO_2$ emissions modeled); energy chain from extraction to final energy; based on SG/MESSAGE example; gas pipe from Siberia+LNG ports. | Discounted energy systems cost minimization as in MESSAGE (LP); trade in gas, coal, electricity and coal; LNG inputs from SEA, MEA, Australia and Alaska, pipelines into China and from Kazakhstan, Irkutsk, Yakutsk, into Japan from Sachalin; exogenous energy use scaled with population; Motivation: to be presented at World Gas Conference in 2003 in Tokyo; plan to introduce economies of scale. |

Table 2-3 Overview of the Literature: Gas Trade Models based on Optimization of Discounted Costs.

Table 2-1, Table 2-2, and Table 2-3, summarize the main characteristics of recent gas trade models. There are basically three dominant classes of gas trade models: macroeconomic models, game theoretic models, and optimization models. Note that the current modeling framework at the ECS Project at IIASA includes both macroeconomic world models (MACRO, MERGE) and a large-scale model (MESSAGE) based on optimization of discounted total costs.

## 3. Players in the simplified version of the Gas Game

| Level | Players | Comments |
|---|---|---|
| **Supply** | Russia | Russian, West Siberian and Caspian gas fields only. The huge East Siberian gas amounts will only be treated in a future extension of the model |
| | "Central Asian Producers" (CAP) | Turkmenistan, Uzbekistan, Kazakhstan |
| | Iran | |
| | LNG from World market | e.g., Algeria, Egypt, Libya, and maybe Nigeria |
| **Demand** | Turkey | (and maybe allow gas piped through Turkey to Western and Eastern Europe) |
| | indigenous demand in producer countries | |
| **Transit Countries** | Azerbaidjan, Armenia, Georgia, Iran, Russia | (later Ukraine) |

Table 3-4 Overview of Players in the Model

The purpose of an initial, simplified version of the model is to model the investments in gas pipelines to Turkey as a dynamic game so as to gain insights in the conditions that determine the time path and the choice of the investments in gas pipelines in the region. The simplified version of the model focuses only on Turkey, whereas a future version will consider extensions to alternative demand markets such as China, India, Japan, and Western Europe (This extended version might eventually explore a gigantic future "gas belt" from Ireland to Australia.) However, we might already allow in the simplified version that some of the gas piped through Turkey will be exported to Western and Eastern Europe.

Key Characteristics for our new Gas Trade Model:
- Non-cooperative game of two-stage investments with profit maximization over pipeline lifetimes (*not* one optimization over the whole time horizon as in MESSAGE)
- price-driven
- time horizon 1998-2050

- competition between long-term contracts and short-term market segments -> analyze gas market liberalization
- analyze first-mover advantage
- Later: key linkage between Asian and European gas markets

## 4. Gas Supply

Fig. 4-1 is an excerpt from the USGS "Gas Futures" map (see [Masters, Turner, 1998]) which illustrates the geographical location of gas fields in the Caspian Sea Region. The color code indicates the potential size of the gas amounts. Note the gigantic amounts in Russia and Iran (Iran holds 15% of the world's gas reserves (http://www.eia.doe.gov/emeu/cabs/caspfull.html)) which are depicted with the color red. Amounts in Turkmenistan and Kazakhstan are at least an order of magnitude smaller (blue color). They are of similar magnitude as North Sea Gas Reserves.



Fig. 4-1. "Gas Futures" map. Potential amounts and locations of gas fields are indicated with different colors. Purple: <0.6, Yellow: from 0.6 to 6, Brown: from 6 to 60, Orange: from 60 to 120, Blue: from 120 to 600, Red: > 600 tcf. Source: U.S. Geological Survey (USGS), see [Masters, Turner, 1998].

Many energy-systems models (e.g., MESSAGE) use cumulative gas supply cost curves. These curves specify the expected cost of gas extraction once a certain cumulative amount of gas has been extracted. In our study we follow (just as in MESSAGE) the classification method of [Rogner, 1997]. He distinguishes 8 different cost categories (see Table 4-1) for gas extraction

according to different economic and technical feasibility. Technological progress in gas extraction continuously transforms higher categories in lower ones.

**Production cost estimates in**

| Categor Source | Natural reserve | Natural EGR | Natural undisc. | Coalbe methan | Tight gas | Gas hydrate | Aquifer gas | Additiona occurrence |
|---|---|---|---|---|---|---|---|---|
| Bourrelier et al. | | | 4 - 30 | > 30 | > 30 | > 30 | > 30 | |
| Dahl and Gjelsvik 1993 | 2 - 15 | | | | | | | |
| EEM Consult 1995 | 7,5 - 14 | | | | | | | |
| Attanasi 1995 | | | 6 - 30 | | | | | |
| Davidson | | | | > 10 | | | | |
| Rogner 1997 | < 10 - | 25 - 29 | 16 - 25 | 29 - 42 | 29 - 42 | d/ | d/ | 42 - 145 |
| Attanasi & Schmoker | | | | > 11 | > 11 | | | |
| Collet and Kuuskraa | | | | | | > 25 | | |
| Category | I | II | III | IV | V | VI | VII | |
| Selected range (1998 | 2 - 16 | 25 - 29 | 4 - 30 | 10 - 42 | 11 - 42 | 25 - 145 | 30 - 145 | |

a/ includes Algeria, FSU and Western Europe (Netherlands, U.K. and

b/ includes Algeria, Iran, Nigeria, Russia and

c/ includes undiscovered resources in onshore and offshore areas of the

d/ gas hydrates and aquifer gas is included in additional occurrences with a cost range of 42 -

e/ quoted in BGR

Table 4-1 Assumed extraction cost ranges for the 8 categories suggested by Rogner. Note that these categories are identical to the current implementation in the MESSAGE model and that overlap in extraction costs are enormous.

In the next step we assume probability distributions for each of the cost categories. For our purpose we use the same regional distributions as assumed in the forthcoming World Energy Assessment report (WEA, 2000). Doing this for natural gas categories I-VI yields cumulative gas supply functions, see Fig. 4-2.



Fig. 4-2. Cumulative gas supply cost curves (Natural Gas Cat. I-VI, based on Rogner method and regional data, adapted by Roehrl).

Finally, note that as specified in the accompanying Excel data file, domestic production in Turkey is assumed to be constant at a level of 0.2 Bcm per year (EIA, 2000).


## 5. Gas Demand

Fig. 5-1 shows natural gas use and production in the 1990s for the players of our Caspian gas game. The difference between production and use is the net exports. Note the strong dominance of Russia both in gas use and exports. Taking this into account and the much larger gas resource base in Russia than in Caspian region countries, it is clear that Central Asian countries will only play an intermediate (a few decades), strategic role (energy security etc.) in Eurasian gas markets. The gas market will clearly be dominated by Russia, and maybe Iran (provided isolation policies of the West are stopped).



Fig. 5-1. Natural gas use and production in the 1990s. The difference between production and use is the net exports.


For our purpose we assume that gas demand is a simple function of GDP and the price of gas:

$D(t) = \alpha P(t)^{e_g} Y(t)^{e_y}$ where D(t) = demand, $\alpha$ = constant, P = price of natural gas, Y = GDP, $e_g$ = price elasticity and $e_y$ is the income elasticity.

Typical price elasticities are around –0.7 for industry in Western-Europe and aggregated over all sectors and markets in 6 countries around –0.93 (see [Golombek et al., 1995]). Brekke et al. (1987) used a price elasticity of –0.7 and an income elasticity of 0.8 for the European continent. It turns out that the respective elasticities for Turkey are quite different (based on 1990s data)[5]: $e_g$ =0.12 and $e_y$ =2.93.

In other words, gas demand is driven mainly by GDP increase and practically insensitive to price changes. That is not surprising for an emerging country with very weak infrastructure and rapidly increasing demand and industrial structure. And as in Asia[6], the fastest growth of gas

---

[5] Note that an average gas price for Turkey in 1998 was about 170 (1995)US$/1000 cm.

[6] In Europe and Asia by far the fastest growth of gas use is expected in the power sector. For example, in Europe more than 50 percent of projected rise in demand for natural gas is due to power sector (see [Capros, 1998]).

use in Turkey is expected in the power sector (which scales perfectly with GDP for both developed and developing countries; see, e.g., [Roehrl, 2000]). Furthermore, the power sector provides the large demand points needed as an "anchor" for new gas infrastructures which are subject to strong scale economies. Also, the gas-combined-cycle technology has emerged as a very competitive ad flexible new option (due to learning by doing since the 60s). Note that Turkish growth in electricity use over the next 15 years is expected to reach at least 8% per year. Currently, 4.2 GW of gas-fired capacity is seeking investment (mainly according to the Turkish "Build-Operate-Transfer"-Model).

However, once a more mature gas market (often characterized by market shares for gas in primary energy of more than 20-25%) is developed, a more even diffusion of gas demand to other sectors is expected. Taking this and the most likely liberalization of energy markets into account, a convergence to elasticities as observed today in Western Europe is to be expected.

Demand for natural gas in Turkey is projected to quadruple within 20 years, and several independent estimates see Turkey's annual consumption reaching as high as 1.4 Tcf by 2020. Gas consumption in Turkey rose from 1.1 Bcm in 1988 to 9.9 Bcm (billion cubic meter) in 1998 (BPAmoco, 1999). In 1998, 6.8 Bcm were imported from Russia by pipeline. 3.0 Bcm was imported from Algeria in the form of LNG. Domestic production was around 0.2 Bcm in 1999 (EIA, 2000) or 2.8% of domestic consumption. Gas consumption is projected to increase to 25.2 bcm in 2010 and 165.2 bcm in 2020 (EIA, 2000).



Fig. 5-2. GDP projections for 1998 to 2050 for the "players" described above. The projections are based on the IPCC SRES marker scenario B2 (see [Riahi, Roehrl, 2000]).

Fig. 5-2 shows GDP projections for 1998 to 2050 for the "players" described above, based on the IPCC SRES marker scenario B2 (see [Riahi, Roehrl, 2000]). B2 is a "dynamics-as-usual" scenario which probably follows a central case future (i.e., not an extreme case). Interestingly, although Iran and CAP countries will experience a considerable population increase in the next 50 years, total GDP of Russia will stay considerably above that of these countries over the whole time horizon. The same will most likely be true for gas demand.

## 6. Gas Pipelines

## 6.1. Overview

The central question here is, of course, how many natural gas pipelines will be needed to transport gas from Central Asia to the Turkish (and European) market?!

Exporters of gas currently have basically 2 options (EIA, 2000): Exporting through the Russian gas pipeline system was the only option available for Caspian gas until 1997. Although over 2 Tcf of Caspian Sea Region gas had been exported via this system in 1990, exports fell to 0.3 Tcf in 1997 because of disputes between Turkmenistan and the Russian gas company Gazprom, a competitor with Turkmenistan, which owned the pipelines. Turkmenistan and Gazprom have come to an agreement to allow Turkmenistan to resume gas exports to Ukraine in 1999. Turkmenistan developed an alternative export route by building a new pipeline from Ekarem (Turkmenistan) to the Iranian border. Limited exports began in 1997, and the ultimate capacity of the pipeline will be 0.5 Tcf (see Section **0**).

Neither of these pipeline options will allow gas from the Caspian Sea Region to compete for a share of the Turkish gas market. Demand for natural gas in Turkey is projected to quadruple within 20 years, and several independent estimates see Turkey's annual consumption reaching as high as 1.4 Tcf by 2020. In addition, there are no export options available to supply gas to the Asian market, where energy demand is expected to grow more rapidly than in any other part of the world.

For natural gas exports to reach these new markets, and for exports to realize their annual potential of 3 Tcf by 2010 and over 5 Tcf within 20 years, several new pipelines will need to be built. Most of the proposals call for pipelines with a capacity of 1 Tcf each, so that 3 or more pipelines could be built, depending upon the extent to which the existing Russian system will be used to supply gas to European markets. Of the proposed new pipelines, the ones that are in the most advanced planning stages are the ones to bring gas to Turkey, such as the Trans-Caspian pipelines and Trans-Iranian lines that are competing to transport 1 Tcf annually to Turkey and European markets (see Section **0** for details).

## 6.2. Current natural gas trade volumes

| Piped gas trade in Eurasia in 1998, in billion cubic meters [bcm] | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | From | | | | | |
| | DK | D | NL | N | UK | Russia | Turkmenistan | Oman | Algeria | **Total imports** |
| **To** | | | | | | | | | | |
| **Europe** | | | | | | | | | | |
| Austria | - | 0.3 | - | 0.4 | - | 5.5 | - | - | - | **6.2** |
| Belgium | - | 0.5 | 5.3 | 5.1 | - | - | - | - | - | **10.9** |

| | | | | | | | | | | Total |
|---|---|---|---|---|---|---|---|---|---|---|
| Bulgaria | - | - | - | - | - | 3.8 | - | - | - | **3.8** |
| Croatia | - | - | - | - | - | 1.1 | - | - | - | **1.1** |
| Czech Republic | - | - | - | 0.8 | - | 8.6 | - | - | - | **9.4** |
| Finland | - | - | - | - | - | 4.2 | - | - | - | **4.2** |
| France | - | - | 5.5 | 10.2 | - | 10.2 | - | - | - | **25.9** |
| Germany | 1.8 | - | 21.1 | 17.5 | 0.9 | 32.3 | - | - | - | **73.6** |
| Greece | - | - | - | - | - | 0.9 | - | - | - | **0.9** |
| Hungary | - | 1.0 | - | - | - | 8.5 | - | - | - | **9.5** |
| Ireland | - | - | - | - | 0.9 | - | - | - | - | **0.9** |
| Italy | - | - | 3.0 | - | - | 16.7 | - | - | 20.9 | **40.6** |
| Luxembourg | - | - | 0.8 | - | - | - | - | - | - | **0.8** |
| Netherlands | - | - | - | 5.2 | 0.6 | - | - | - | - | **5.8** |
| Poland | - | - | - | - | - | 7.5 | - | - | - | **7.5** |
| Portugal | - | - | - | - | - | - | - | - | 0.9 | **0.9** |
| Romania | - | - | - | - | - | 3.8 | - | - | - | **3.8** |
| Slovakia | - | - | - | - | - | 6.9 | - | - | - | **6.9** |
| Slovenia | - | - | - | - | - | 0.5 | - | - | 0.4 | **0.9** |
| Spain | - | - | - | 2.5 | - | - | - | - | 4.5 | **7.0** |
| Sweden | 0.9 | - | - | - | - | - | - | - | - | **0.9** |
| Switzerland | - | 1.5 | 0.7 | - | - | 0.5 | - | - | - | **2.7** |
| Turkey | - | - | - | - | - | 6.8 | - | - | - | **6.8** |
| UK | - | - | - | 0.9 | - | | - | - | - | **0.9** |
| Others | - | - | - | - | - | 2.5 | - | - | - | **2.5** |
| **Middle East** | | | | | | | | | | |
| Iran | - | - | - | - | - | - | 1.8 | - | - | **1.8** |
| United Arab Emirates | - | - | - | - | - | - | - | 0.5 | - | **0.5** |
| **Africa** | | | | | | | | | | |
| Tunisia | - | - | - | - | - | - | - | - | 0.8 | **0.8** |
| **TOTAL EXPORTS** | **2.7** | **3.3** | **36.4** | **42.6** | **2.4** | **120.3** | **1.8** | **0.5** | **27.5** | **237.3** |

Table 6-1 Natural Gas Trade through pipelines in Eurasia in 1998. Data source: CEDIGAZ.

| World LNG Trade in 1998, in billion cubic meters [bcm] | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | From | | | | | | | | |
| | USA | Qatar | UAE | Algeria | Libya | Australia | Brunei | Indonesia | Malaysia | **Total imports** |
| **To** | | | | | | | | | |
| **North America** | | | | | | | | | |
| USA | - | - | 0.1 | 2.0 | - | 0.2 | - | - | - | **2.3** |
| **Europe** | | | | | | | | | |
| Belgium | - | - | - | 4.3 | - | - | - | - | - | **4.3** |
| France | - | - | - | 9.8 | - | - | - | - | - | **9.8** |
| Italy | - | - | 0.1 | 1.9 | - | - | - | - | - | **2.0** |
| Spain | - | 0.5 | 0.6 | 3.9 | 0.9 | - | - | - | - | **5.9** |
| Turkey | - | 0.6 | - | 3.0 | - | - | - | - | - | **3.6** |
| **Asia Pacific** | | | | | | | | | |
| Japan | 1.8 | 3.7 | 6.2 | - | - | 9.7 | 7.3 | 24.2 | 13.2 | **66.1** |
| South Korea | - | - | 0.1 | - | - | - | 0.8 | 9.5 | 3.9 | **14.3** |
| Taiwan | - | - | - | - | - | - | - | 2.4 | 2.3 | **4.7** |
| **TOTAL EXPORTS** | **1.8** | **4.8** | **7.1** | **24.9** | **0.9** | **9.9** | **8.1** | **36.1** | **19.4** | **113.0** |

Table 6-2 World LNG Trade in 1998, in billion cubic meters. Source: CEDIGAZ.

## 6.3. Existing and Planned Pipeline Project

Supply to Turkey can come from the Russian Federation, Iran, Turkmenistan (and Kazakhstan, Azerbaijan), as well as Egypt and Algeria. This Section provides a brief overview of the major proposed pipeline projects, see Fig. 6-1.

## Gas Export Projects to Turkey:

a)  Russia-Ukraine-Bulgaria-Turkey

Through this existing pipeline 280 bcf were delivered to Turkey in 1998. In December 1997 another 25 year deal was signed between Gazprom and Turkey to deliver 500 bcf annually by early 2005. To deliver these large amounts either large investments in the existing pipeline route and contracts with transit countries are needed, or new pipeline connections will be built, e.g.:

b) "Blue Stream"

This project proposes a direct connection between Russia and Turkey under the Black Sea (to avoid trouble with transit countries). This 758 miles long connection would be the deepest underwater gas pipeline in the world. Therefore it is technically challenging, expensive (4 US$ billion), and probably won't be finished before 2003.

a) Turkmenistan into Iran (already exists)

This short connection (only about 190 mill. US$ [7]) is already existing. Gas exports to Iran began in late 1997 (141 bcf). A capacity extension to 283 bcf is planned for 2006, together with an extension to Teheran.

b) SWAP deal between Turkmenistan, Iran and Turkey

Example c) may become part of a large-scale SWAP deal between Turkmenistan, Iran and Turkey. This plan would require three new pipelines to Turkey (46 inch). Although economically reasonable, it faces fierce US opposition, which is targeted to continue to isolate Iran in the region. Supporters of the SWAP deal however claim that it would not violate the American Iran-Libya Sanctions Act (ILSA[8]), because Iran would only receive transit fees for moving gas to Turkey, rather than exporting gas themselves.

c) Agreement between Turkey and Iran

In 1996, Iran and Turkey signed a $20-billion agreement that calls for Iran to supply Turkey with natural gas over a period of 22 years. Exports of Iranian gas to Turkey were slated to start in 1999 at an initial rate of 300 Mmcf/d and rise to a level of 1,000 Mmcf/d in 2005. In November 1998, Turkey began construction of a 623-mile pipeline that could transport gas westward from Iran. In January 2000, Iran said that it accepted Turkey's request to delay the purchase of Iranian natural gas until September 2001. Turkey said that it had been unable to complete its portion of the pipeline due to economic problems. With respect to Iran that part of the pipeline that runs on the Iranian side has been completed (see [Ignatius, 2000]). Construction on the Turkish side of the 160 mile, line from the Iranian border to Erzurum has been slowed.

d) US to priority: "Trans-Caspian" Gas Pipeline (TCGP)

For obvious geo-political considerations (to isolate Russia and Iran) the US government pushes the TCGP project. This would be a 1050 miles gas pipeline from Turkmenistan underneath the Caspian Sea to Baku, through Azerbaijan and Georgia onto Turkey for approximately 2.5 billion US$. A first agreement was signed in May 1999 to deliver 565 bcf gas by as early as 2002. These large amounts are intended to be exported further to Europe in the future. In parallel with the US ideas, this route is also a Turkish priority. However, recently a large new natural gas field was found in Azerbaijan, which makes the potential transit country for Turkmen gas to become a direct competitor. Note also that Iran strongly objects to this gas route (on "environmental" grounds) and exerts pressure on Turkmenistan.

---

[7] Note that for the model implementation, annual variable operation and maintenance (vom) costs have to be added. They are around 5-10% of the pipeline investment costs. Furthermore, a plant life of around 30 years is assumed for gas pipelines.

[8] The US act "ILSA" was signed in 1996 and imposes mandatory and discretionary sanctions on non-US companies which invest more than 20 million US$ annually in the Iranian oil or gas sector.

e) North Caucasus-Transcaspian Gas Pipeline (NCTGP)

The NCTGP provides a link round the Caspian-Russia into Georgia-Armenia-Turkey. Deliveries through this line were cut off in 1997 due to lack of payments for gas deliveries by Georgia. Furthermore, the link from Georgia to Armenia was destroyed in the civil war (1995). Only 250 million US$ would be needed to upgrade the pipelines and carry 425 bcf/year to Turkey. Due to American pressure Gazprom's partner Royal Dutch/Shell has switched recently from this project to the TCGP. Consequently, Gazprom now favors the much more expensive "Blue Stream".

f) LNG ports planned in Kazakhstan, Turkmenistan, Baku and Georgia

With the intention to export LNG through the Black Sea to Europe, a number of LNG ports are planned in the region (planned investment of about 250 million US$).

g) MOU signed between Turkey and Egypt

In June 1998 a MOU was signed between Turkey and Egypt to build a
- Offshore pipeline under the Mediterranean to deliver Egyptian gas to Gaza, Israel, Egypt, Lebanon, Syria, and South-East Turkey
- LNG ports to deliver Egyptian LNG to Izmir. This would require a 1.2 billion US$ liquefaction terminal.


Apart from these projects there are a number of other important projects which might be dealt with in a future model extension, e.g.,

- China pipeline: 5000 miles from Turkmenistan to China (with a possible extension to Japan), 1 tcf/yr, 8.5 billion US$
- Centgas: Pipeline from Turkmenistan via Afganistan to Pakistan (and probably further on to India); 900 miles, 2-2.5 billion US$.
- Pipelines[9] down to Southern-Iran with subsequent LNG terminals to export LNG to Asian countries (including Japan).
- Pipeline Iran-Pakistan-India, where a MOU was signed in Spring 2000.

---

[9] Although domestic gas consumption is growing rapidly, including use as a motor fuel, Iran continues to promote export markets for its natural gas. Possibilities include pipelines to Turkey, Armenia, Europe, Pakistan, and India, plus the possibility of an LNG facility for producing exports to Asia.

Fig. 6-1. Gas Pipeline Routes to Turkey (Source: Gas Matters).

A (slightly outdated, but otherwise superb) overview of the pipelines projects (both oil and gas) in the region can be found on http://www.eia.doe.gov/emeu/cabs/caspfull.html (Source: EIA, 2000).

Finally, Fig. 6-2 provides an overview of the main existing pipeline network to deliver gas from the countries of the Former Soviet Union (mainly Russia, but previously also Kazakhstan and Turkmenistan) to Western and Eastern Europe. Note the large number of transit countries, which are responsible for much of Gazprom's deficit due to non-payments issues.



Fig. 6-2. Main export pipelines from the countries of the Former Soviet Union (FSU) to Western and Eastern Europe. Source: Gas Matters.

From the discussion above it seems evident that geo-political consideration seem to be as important as economic and financial considerations for pipeline projects to tap the large resources in Central Asia. Other political risks include also uncertainties because of the non-payments issue for gas deliveries, as well as unresolved Caspian Sea legal issues. Is the Caspian an inland lake or governed by the Law of the Sea? The legal situation of an inland lake would support joint development projects. However, if the Caspian is governed by the law of the sea full maritime boundaries would apply.

## 7.    The Model of Optimal Investment in Gas Pipelines

Competition process on markets of gas pipeline construction and commercialization is formalized as a dynamic nonzero sum game. Dynamics of a competitor describes investment process of pipeline construction with discount, obsolescence and delay effects. It is assumed that a competitor controls the system of investment exponential trajectories by selecting the time of pipeline commercialization, balancing the level of gas supply and optimizing the level of current investment. Revenues of gas sales for one gas supply project depends on dynamics of gas prices (and elasticity of gas prices) which in turn are determined essentially by growing demand on a gas market and by consequences of commercialization times and gas supply of other projects. The model takes into account dynamics of costs of gas production at different gas fields and costs of gas transportation. A competitor can make a decision on the commercialization time and gas supply of its project dynamically taking into account information about the current stage and dynamics of other competing gas projects. We propose a version of solution when a competitor decomposes the control process into four dynamically interacting stages: assessment of potential gas demand on markets and prediction of commercialization times and gas supply of active gas projects; regulation of gas supply through active pipelines; scenario selection for the commercialization time of the controlled gas projects; optimization of investment level into constructing pipelines.

### 7.1   Objects

Let $m = 1, \cdots, M$  be the number of gas markets;

$j = 1, \cdots, J(m)$ - the number of consumers on market $m$ ;

$l = 1, \cdots, L$ - the number of gas fields;

$i = 1, \cdots, I(l,m)$ - the number of pipelines from field $l$ to market $m$ ;

### 7.2   Supply Part

Construction and gas supply of pipeline $i$ is described by the following parameters:

$t = t_0 = t_i^0$ - the initial time of construction of pipelines;

$t_i^a$ - the time of gas supply by pipeline $i$ ;

$s \in [t_0, t_i^a]$ - the current time;

$x_i(s)$ - the accumulated investment at time $s$ in pipeline $i$ ;

$x_i^a = x_i(t_i^a)$ - the level of accumulated investments which is necessary for starting gas supply;

$r_i(s)$ - the current investment into pipeline $i$ ;

$y_i^m$ - gas supply of pipeline $i$  to market $m$ starting from time $t_i^a$ ;

$0 < \gamma < 1, \alpha = 1/\gamma > 1$ - delay coefficient of investment;

$\lambda > 0$ - discount of investment;

$\sigma > 0$ - obsolescence coefficient of investment;

$k = 1, \cdots, K_m, K^m = K^m(\tau)$, the quantity of gas pipelines to market $m$ from gas fields;

$y^m = y^m(\tau)$ - total gas supply to market $m$ at time $\tau$

$$y^m = y^m(\tau) = \sum_{k=1}^{K^m} y_k^m(\tau), \qquad 0 \le y_k^m(\tau) \le \overline{y}_k^m \tag{1}$$

Here $\overline{y}_k^m$ is the exogenous maximum supply of pipeline $k$ to market $m$

## 7.3 Dynamics of Investment into Pipeline

The dynamics of the accumulative investment $x_i$ is presented by the following differential equation

$$\dot{x}_i(s) = -\sigma x_i(s) + r_i^\gamma(s), \qquad x_i(t_0) = x_i^0, \qquad x_i(t_i^a) = x_i^a \tag{2}$$

with the obsolescence parameter $\sigma$ of the accumulative investments $x_i$ and the delay coefficient $\gamma$ of the present-day investments $r_i$.

## 7.4 Cost of Pipeline Construction

Cost of pipeline construction is measured by minimum expenditures $W_i$

$$\int_{t_0}^{t_i^a} e^{-\lambda s} r_i(s) ds \to \min_{r_i(\cdot)} \to W_i = W_i(t, x_i, t_i^a, x_i^a, \gamma, \lambda, \sigma) \tag{3}$$

Here minimum expenditures are expressed by solution of the optimal control problem

$$W_i = \rho^{(\alpha-1)} \frac{e^{-\lambda t_i^a}(x_i^a - x_i e^{-\sigma(t_i^a - t)})^\alpha}{(1 - e^{-\rho(t_i^a - t)})^{(\alpha-1)}}, \rho = \frac{(\alpha\sigma + \lambda)}{(\alpha - 1)} \tag{4}$$

The corresponding optimal feedback is described by relations

$$r_i^\gamma = r_i^\gamma(t, x_i, t_i^a, x_i^a, \alpha, \lambda, \sigma) = \rho \frac{(e^{\sigma(t_i^a - t)} x_i^a - x_i)}{(e^{\rho(t_i^a - t)} - 1)} \tag{5}$$

Optimal innovation trajectory is given by formula

$$x_i(s) = e^{-\sigma(s-t)}(x_i + \frac{(e^{\sigma(t_i^a - t)} x_i^a - x_i)(e^{\rho(s-t)} - 1)}{(e^{\rho(t_i^a - t)} - 1)})$$

## 7.5 Demand Part

The demand part is determined by the following variables:

$Y_j$ - production function of consumer (country) $j$ ;

$e_y$ - elasticity of production $Y_j$ ;

$Q_j$ - the indigenous gas supply of consumer $j$ .

The demand of consumer $j$ on market $m$ is given by function of Cobb-Douglas type

$$d_j^m = P_g^{-e_g} (A_j Y_j^{e_g} - Q_j) \tag{6}$$

The total demand on market $m$ is described by relation

$$d^m = \sum_{j=1}^{J_m} d_j^m = P_g^{-e_g} \sum_{j=1}^{J_m} (A_j Y_j^{e_y} - Q_j) \tag{7}$$

## 7.6 Price Formation Mechanism

The price on gas is defined according to the principle of stabilizing demand $d^m$ and supply $y^m$ on the market $m$ by the price mechanism

$$d^m = P_g^{-e_g} z^m \Leftrightarrow P_g = P_g(s) = (\frac{z^m}{y^m})^{1/e_g} \tag{8}$$

Here symbol $z^m$ stands for the potential gas consumption on market $m$

$$z^m = \sum_{j=1}^{J_m} (A_j Y_j^{e_y} - Q_j) \tag{9}$$

## 7.7 Revenues of Gas Sales

Revenues of gas sales of pipeline $i$ is described by integrated consumption of its supply $y_i = y_i(s)$ at a market price $P_g = P_g(s)$

$$V_i = V_i(t_i^a, y_i, y^m, z^m) = \int_{t_i^a}^{+\infty} e^{-\lambda s} P_g(s) y_i(s) ds = \int_{t_i^a}^{+\infty} e^{-\lambda s} (\frac{z^m(s)}{y^m(s)})^{1/e_g} y_i(s) ds \tag{10}$$

## 7.8 Cost of Gas Production at Field $l$ for Pipeline $i$

The cost $C_i^l$ of gas extraction at field $l$ for pipeline $i$ is given by relation

$$C_i^l = C_i^l(t_i^a) = \int_{t_i^a}^{+\infty} e^{-\lambda s} P^l y_i(s) ds \tag{11}$$

Here $P^l$ is the price of the gas production at field $l$. It may depend on cumulative extraction $C^l$ at field $l$

$$P^l = P^l(C^l), C^l = C^l(s) = \int_t^s (\sum_{i=1}^{I_t} y_i^l(\tau)) d\tau$$

Here $I_t$ is the number of pipelines at the field $l$ and $C \to P^l(C)$ is a monotonically increasing function.

## 7.9 Cost of Gas Transportation from Field $l$ to market $m$ through Pipeline $i$

The cost $C_i^{l,m}$ of gas transportation from field $l$ to market $m$ depends on prices $P_i^n$ of gas transit through countries $n = 1, \cdots, N_i$

$$C_i^{l,m} = C_i^{l,m}(t_i^a) = \int_{t_i^a}^{+\infty} e^{-\lambda s} (\sum_{n=1}^{N_i} P_i^n) y_i(s) ds \tag{12}$$

## 7.10 Profit of Pipelines: Revenues Minus Costs

The profit of the exploited pipeline $i$ from field $l$ to market $m$ is defined as usual: revenues minus costs

$$G_i = G_i^{l,m}(t, x_i, t_i^a, x_i^a, y_i, y^m, z^m) =$$

$$V_i(t_i^a, y_i, y^m, z^m) - W_i(t, x_i, t_i^a, x_i^a) - C_i^l(t_i^a) - C_i^{l,m}(t_i^a) =$$

$$\int_{t_i^a}^{+\infty} e^{-\lambda s} \left(\frac{z^m(s)}{y^m(s)}\right)^{1/e_g} y_i(s) ds - W_i(t, x_i, t_i^a, x_i^a) -$$

$$\int_{t_i^a}^{+\infty} e^{-\lambda s} P^l y_i(s) ds - \int_{t_i^a}^{+\infty} e^{-\lambda s} (\sum_{n=1}^{N_t} P_i^n) y_i(s) ds \tag{13}$$

## 7.11 Equilibrium in Gas Price and Supply

The prices of the unbound market at the time $t$ are presented by the relation

$$P(t, y_1, ..., y_i, ..., y_K) = \frac{(z^m(t))^{1/e_g}}{(y_1 + ... + y_i + ... y_K)^{1/e_g}} \tag{14}$$

The prices of the energy resources equivalent to gas, i.e. the prices on the liquid natural gas (LNG), are given as the exogenous value at the time $t$

$$P_{LNG} = P_{LNG}(t) \tag{15}$$

Profit of the player $i$ under the market price $P$ is determined by the formula

$$R_i^{(1)}(t, y_1, ..., y_i, ..., y_K) = P(t, y_1, ..., y_i, ..., y_K) y_i - P_i y_i \tag{16}$$

125

$$P_i = P^l + \sum_{n=1}^{N_t} P_i^n$$

Profit of the player $i$ under the LNG price $P_{LNG}$ is calculated by the following relation

$$R_i^{(2)}(t, y_1, ..., y_i, ..., y_K) = P_{LNG}(t) y_i - P_i y_i \qquad (17)$$

The real profit of the player $i$ is estimated by the function of the minimum type

$$R_i(t, y_1, ..., y_i, ..., y_K) = \min\{R_i^{(1)}(t, y_1, ..., y_i, ..., y_K), R_i^{(2)}(t, y_1, ..., y_i, ..., y_K)\}$$

A Nash equilibrium $(y_1^*, ..., y_i^*, ..., y_K^*)$ in the space of supply parameters is defined by relations

$$R_i(t, y_1^*, ..., y_i^*, ..., y_K^*) = \max_{0 \le y_i \le \bar{y}_i} R_i(t, y_1^*, ..., y_i, ..., y_K^*), \quad i = 1, ..., K \qquad (18)$$

Under conditions of concavity of profit functions $y_i \to R_i(t, y_1, ..., y_i, ..., y_K)$ a Nash equilibrium exists according to the Kakutani theorem on fixed points.

The equilibrium gas price $P_g$ is determined as the price calculated at a Nash equilibrium

$$P_g = P_g(t) = \min\{P(t, y_1^*, ..., y_i^*, ..., y_K^*), P_{LNG}(t)\} \qquad (19)$$

## 7.12  Algorithm for Searching Nash Equilibrium

For finding a Nash equilibrium the "best reply" dynamics can be used

$$\dot{y}_i(\tau) = S(-sign(y_i(\tau)) + sign(\bar{y}_i - y_i(\tau)) + sign(\partial_{y_i} R_i(t, y_1(\tau), ..., y_i(\tau), ... y_K(\tau)))) \quad (20)$$

Here parameter $S$ is the regulator rate of the fast time $\tau$. The derivative $\partial_{y_i} R_i$ is determined by the formula

$$\partial_{y_i} R_i(t, y_1, ..., y_i, ... y_K) = \begin{cases} \dfrac{\partial R_i^{(1)}}{\partial y_i}, R_i^{(1)} < R_i^{(2)} \\[2mm] \dfrac{\partial R_i^{(2)}}{\partial y_i}, R_i^{(1)} \ge R_i^2 \end{cases} \qquad (21)$$

Under the certain convergence conditions on parameters of the model the "best reply" dynamics stabilizes at a Nash equilibrium.

## 7.13  Continuous Development of the Market Demand

It is naturally to assume that demand of the market develops continuously depending on inertial dynamics of production $Y_j$

$$(\ln Y_j(s))'' = v_j(s), Y_j(t_0) = Y_j^0, \frac{\dot{Y}_j(t_0)}{Y_j(t_0)} = \delta_j^0, |v_j(s)| \le v^0 \tag{22}$$

According to exponential trajectories

$$Y_j(s) = Y_j^0 e^{\delta_j(s-t_0)(s-t_0)}, \delta_j(s-t_0) = \delta_j^0 + \varepsilon_j(s-t_0), |\varepsilon_j(s-t_0)| \le \varepsilon^0 \tag{23}$$

## 7.14    Discontinuous Character of Gas Supply

Due to the fact that the number of supplying pipelines on the market $m$ depends on time $K^m = K^m(\tau)$ it is clear that the integral function in revenues $R_i$ is piecewise differentiable since it has the piecewise continuous integrands. Calculation of its values can be reduced to a finite series of integrals with continuous integrands.

## 7.15    Clusters of Inertial Suppliers

Assuming that gas suppliers can be grouped into inertial clusters $r = 1,...,R_m$ one can describe dynamics of investments $X_r$ into new pipelines on market $m$ by inertial dynamics

$$(\ln X_r(s))'' = w_r(s), X_r(t_0) = X_r^0, \frac{\dot{X}_r(t_0)}{X_r(t_0)} = -\sigma + \eta_r^0, |w_r(s)| \le w^0 \tag{24}$$

according to exponential trajectories

$$X_r(s) = X_r^0 e^{\eta_r(s-t_0)(s-t_0)}, \eta_r(s-t_0) = -\sigma + \eta_r^0 + \xi_r(s-t_0), |\xi_r(s-t_0)| \le \xi^0 \tag{25}$$

and predict the exploitation time $t_r^a$ of the clustered suppliers $X_r$.

## 7.16    Maximization of Profit

The task of optimization of investments into pipeline $i$ is to maximize profit $G_i$ (13) by selecting the optimal time $t_i^a$ for starting exploitation, minimizing the current level of investments $x_i(s)$ and choosing actual levels $y_i$ of gas.

## 8.    Dynamic Optimization Algorithm

Let us describe the algorithm of optimization of the gas supply, the assessment of the market dynamics, optimal timing of the gas pipeline commercialization and optimal investment into the   construction of a pipeline. The algorithm is constructed on the basis of the feedback principle which allows to react on the changing situation and to fight with uncertainties. The peculiarity of this algorithms consists in its dynamic setting when all four optimization levels of the model dynamically interact with each other, interchange information and decision making rules.

## 8.1 Optimization of Supply

At the current time $t$ the gas supply is stabilized at a Nash equilibrium point $(y_1^*(t),..., y_i^*(t),..., y_K^*(t))$ (18) with an equilibrium gas price $P_g(t)$ (19) and an investor $i$ supplies gas at the level $y_i = y_i^*(t)$.

## 8.2 Assessment of the Market Dynamics

The evaluation of the market GDP growth $(Y_j(t), \dot{Y}_j(t))$ according to dynamics (22)-(23) and the estimation of the development parameters $(X_r(t), \dot{X}_r(t))$ on the basis of equations (24)-(25) give an investor $i$ the opportunity to forecast the future situation on the market: to predict the commercialization times $t_r^a$ of the opponents; to evaluate for times $\tau \geq t_i^a$ after commercialization the value of the potential gas demand $z^m(\tau)$ (9), a Nash equilibrium supply $(y_1,..., y_i,..., y_K) = (y_1^*(\tau),..., y_i^*(\tau),..., y_K^*(\tau))$ (18), the amount of the total gas supply $y^m(\tau)$ (1), and, as a consequence, to project the gas price $P_g(\tau)$ (19).

## 8.3 Selection and Optimization of the Commercialization Time

An investor $i$ can optimize the commercialization time $t_i^a$ for the constructed pipeline, maximizing the profit function $G_i(t, x_i, t_i^a, x_i^a, y_i, y^m, z^m)$ (13) over the time variable $t_i^a$, and fix the current scenario of the gas pipeline construction.

## 8.4 Optimal Synthesis of the Investment Level

An investor $i$ substitutes the optimal commercialization time $t_i^a$ calculated for the maximum profit $G_i$ (13) into the minimum cost function $W_i$ (4) of the gas pipeline construction and selects the optimal investment plan $r_i$ (5).

## 9. Conclusion

The paper is devoted to the problem of an adequate assessment of the future energy infrastructure in energy-consuming regions. As a demonstrative example we consider energy routes in the Caspian region to the energy market in Turkey. Our basic idea is to present analysis as dynamic interaction of the existent huge projects of gas pipeline construction. The main competitors are the Russian, Turkmenian and Iranian gas pipeline projects. In the general model of competition between gas pipelines it is reasonable to distinguish several control

levels of the process regulation: current investment into a new gas pipeline, future supply regulation, selection of commercialization times, and link them through the price formation mechanism. We make an accent on a game model of timing in which the main control parameters of the competitors are times of finalizing of gas pipeline projects. The model captures Nash multiequilibria situation in this competition and provides an instrument for reasonable dynamic selection of an equilibrium scenario for the gas market.

## References

1.  Arrow, K.J., Kurz, M., 1970, Public Investment, the Rate of Return and Optimal Fiscal Policy, Baltimore, Johns Hopkins University Press.

2.  Barzel, Y., 1968, Optimal Timing of Innovations, The Review of Economics and Statistics, Vol. 50, pp. 348-355.

3.  Basar, T., Olsder, G.J., 1982, Dynamic Noncooperative Game Theory, London, N.Y., Acad. Press.

4.  Beltramo, M.A., Manne, A.S., Weyant, J.P., 1986, A North American Gas Trade Model, (GTM), *The Energy Journal*, Vol. 7, No. 3.

5.  Berg, E., Boug P., Kverndokk, S., 1998, Norwegian Gas Sales and the Impacts on the European CO2 Emissions, Nota di Lavoro 9.98, prepared for Fondazione Eni Enrico Mattei.

6.  Boots, M.G., Van Oostvorn, F., 1999, Impacts of Market Liberalization on the EU Gas Industry, The Shared Analysis Project Energy Policy in Europe and Prospects to 2020, Vol. 9, Prepared for the European Commission Directorate General for Energy.

7.  Brekke, K.A., Gjelsvik, E., Vatne, B.H., 1987, A Dynamic Supply Side Game Applied to the European Gas Market, Discussion Paper No. 22, Central Bureau of Statistics, Oslo, Norway.

8.  Brekke, K.A., Gjelsvik, E., Vatne, B.H., 1991, A Dynamic Investment Game - The Fight for Market Shares in the European Gas Market, Manuscipt.

9.  Capros, 1998, Baseline Projections for Shared Analysis, Manuscipt.

11. EIA, 2000, Caspian Sea Region, June, 2000, United States, Energy Information Administration, Washington.( http://www.eia.doe.gov/emeu/cabs/caspfull.html).

12. Fujii, Y., Yamaji, K., 1999, An Energy Infrastructure Model for Asia/Eurasia, presented at Joint Energy Meeting, IEA, Paris, 16-18 June, 1999.

13. Golombek, R., Gfelsvik, E., Rosendahl, K.E., 1995, Effects of Liberalizing the Natural Gas Markets in Western Europe, The Energy Journal, Vol. 16, No. 1, P. 85-111.

14. Hofbauer, J., Sigmund K., 1988, The Theory of Evolution and Dynamic Systems, Cambridge, Cambridge Univ. Press.

15. Ignatius, D., 2000, The Great Game Gets Rough, The Washington Post, January 26, page A23.

16. Intriligator, M., 1971, Mathematical Optimization and Economic Theory, N.Y., Prentice-Hall.

17. Klaassen, G., McDonald, A., Zhao, J., 2001, The Future of Gas Infrastructure in Eurasia, Energy Policy, Vol. 29, P. 399-414.

18. Krasovskii, N.N., Subbotin, A.I., 1988, Game-Theoretical Control Problems, N.Y., Berlin, Springer.

19. Kryazhimskii, A., Nentjes, A., Shibayev, S., Tarasyev, A., 2001, Modeling Market Equilibrium for Transboundary Environmental Problem, Nonlinear Analysis: Theory, Methods and Applications, Vol. 42, P. 991-1002.

20. Makarov, A., 1999, Diversification of Russian Gas Export Routes. Proceedings of International Conference "The Role of Russian and CIS Gas Countries in Deregulated Energy Markets", Paris, The Moscow International Energy Club/Universite Paris Dauphine, Centre de Geopolitique at de l'Energie et des Matieres Premieres, December 6-7, 1999.

21. Masters, C., Turner, R.M., 1998, World Petroleum Futures (Gas), U.S. Geological Survey, Open File Report 98-486.

22. Pontryagin, L.S., Boltyanskii, V.G., Gamkrelidze, R.V., Mishchenko, E.F., 1962, The Mathematical Theory of Optimal Processes, N.Y., Interscience.

23. Riahi, K., Roehrl, A.R., 2000, Energy Technology Strategies for Carbon Dioxide Mitigation and Sustainable Development , Environmental Economics and Policy Studies, Vol. 3, No. 2, P. 89-123.

24. Rogner, H.-H., 1997, An Assessment of World Hydrocarbon Resources, Annual Review of Energy and Environment, Vol. 22, P. 217-262.

25. Roehrl, R.A., 2000, A Spatial Model of Electricity Use in China - Disaggregating Global Energy Model Results Based on Economic Growth Theory, IR-00-012, International Institute for Applied Systems Analysis, Laxenburg.

26. Smeers, Y., Wolf, D.D., 1997, A Stochastic Version of a Stackelberg-Nash-Cournot Equilibrium Model, Management Science, Vol. 43, No. 2.

27. Sugiyama, T., 1999, China Provincial Energy and Emissions Model (CPE Model), Manuscript.

28. Tarasyev, A.M., Watanabe, C., 2001, Dynamic Optimality Principles and Sensitivity Analysis in Models of Economic Growth, Nonlinear Analysis: Theory, Methods and Applications, Vol. 47, No. 4, P. 2309-2320.

**Units**

1 bcm = 35.31466672 bcf

bcm: billion cubic meters

bcf: billion cubic feet

tcf: trillion cubic feet

**Appendix: Overview Tables (other sources)**

| Russian Federation | | |
| --- | --- | --- |
| Existing export (1998) | 6.8 bcm | |
| Excess Capacity | (7.8-6.8) bcm ? | |
| Investment Projects | Capacity | Investment |
| Blue Stream | 14 bcm | US$4 billion (1999) (758 miles) |
| Bulgarian route | ? | |
| Armenian route | ? | |

Table A-1 The Russian Federation. Source: EIA (2000), BPAmoco (1999).

.

| Turkmenistan | | |
|---|---|---|
| Existing export (1998) | 0 bcm | |
| Excess Capacity | 0 | |
| Investment Projects | Capacity | Investment |
| Transcaspian | 15.8 bcm | US$2.5 billion (1999) (1050 miles) |
| Iranian Swap route | ? | |

Table A-2 Turkmenistan. Source: EIA (2000), BPAmoco (1999).

| Iran | | |
|---|---|---|
| Existing export (1998) | 0 bcm | |
| Excess Capacity | 0 | |
| Investment Projects | Capacity | Investment |
| Pm | 15.8 bcm | US$2.5 billion (1999) (1050 miles) |
| Iranian Swap route | ? | |

Table A-3 Iran. Source: EIA (2000), BPAmoco (1999).

# Optimal control of dynamic system in the presence of risky factors

*Sergei Aseev*

## Abstract

This paper deals with problems of optimal control for dynamical systems whose state spaces contain domains of risk. Serious difficulties in the analysis of such problems arise due to the discontinuities with respect to the state variable, which may occur in the system's dynamics or in the cost functional. Two problems of optimal control with domains of risk are considered: a problem with state constraints and a problem of time-optimal crossing a given domain.

# 1. Introduction.

The classical optimal control problem *(P)* is formulated as follows:

$$\dot{x} = f(x,u), \quad u \in U;$$
$$x(0) = x_0, \quad x(T) = x_1;$$
$$J(x(\cdot),u(\cdot)) = \int_0^T f^0(x(t),u(t))\,dt \to \min.$$

Here $x = (x^1,\ldots,x^n) \in R^n$ is a state vector, $u = (u^1,\ldots,u^m) \in R^m$ is a control parameter, $U$ is a nonempty compact subset of $R^m$, a time interval $[0,T]$ is free and both initial state $x_0$ and final state $x_1$ of the system are fixed. The set of admissible controls consists of all bounded measurable vector functions $u(t)$ such that $u(t) \in U$ almost everywhere (a.e.) on $[0,T]$.

The aim of the process of control consists in steering the system from the initial state $x_0$ to the final state $x_1$ in such a way that the cost functional $J(x(\cdot),u(\cdot))$ takes its minimal value.

The formulated optimal control problem *(P)* is a classical one; the corresponding theory is well developed [12].

An important feature of this problem consists in the fact that both vector function $f(x,u)$ and scalar function $f^0(x,u)$ are assumed to be smooth in respect to state variables. This smoothness assumption plays an important role in the basic constructions of the classical optimal control theory [12]. Note that recently a corresponding theory was constructed also for optimal control problems for which the above functions are only Lipschitzian in respect to the state variables [9].

The present paper is concerned with a more complicated situation when the setting of an optimal control problem admits discontinuities in the state variables. Such situation naturally occurs in the case when the statement of the problem involves a given *risky* domain $Z$ in the state space $R^n$. In this case the supplement $G = R^n \setminus Z$ of the set $Z$ may be considered as a *safety* domain for the system. The term *risky* domain means that the situation when the system's trajectory $x(t)$ hits set $Z$ is possible but undesirable by safety reasons. The optimal control problems discussed in the present paper consist of finding an admissible control $u(t)$ such that it brings the system from the initial state $x_0$ to the final state $x_1$ by such a way that the cost functional $J(x(\cdot),u(\cdot))$ takes its minimal value and, simultaneously, the time interval, during which the system's trajectory stays in the risky domain $Z$, is minimal.

The optimal control problems of this type arise in various processes in technology, economics and ecology. These optimal control problems have many specific features and can be formalized differently. The deterministic formalizations of such problems often involve the discontinuities with respect to the state variables in the system's dynamics or in the cost functional.

In this paper two qualitatively different cases are considered.

The first case is related to the situation when there exists at least one admissible control, which transfers the system from initial state $x_0$ to the final state $x_1$ and such that the corresponding trajectory $x(t)$ of the controlled system stays away from the risky domain $Z$. In this case the safety domain $G$ can be selected as a natural state constraint and the corresponding optimal control

problem with state constraints may be considered as an adequate mathematical model for the control process under consideration.

The second case is characterized by a situation when there is no admissible control which steers the system from $x_0$ to $x_1$ so as the corresponding trajectory $x(t)$ stays away from the set $Z$. It may happen, for example, if the initial state $x_0$ or the final state $x_1$ belongs to the risky domain $Z$.

In the present paper both possible situation are considered. The main attention is paid to the discussing recently developed necessary optimality conditions for the corresponding problems.

## 2. Optimal control problem with state constraints.

The state constraints naturally arise in the analysis of many optimal control problems. The origin of these constraints may be essentially diverse.

On the one hand, the state constraints can be of a physical nature. For example, in the analysis of the control of an airplane, one can consider the Earth's surface as a natural state constraint. In this case, the violation of the state boundary is impossible, and a control system (a plane) interacts with the state boundary during landing; this interaction must be taken into account in the mathematical model of the control processes.

On the other hand, the state constraints may arise as artificial constraints imposed in view of certain external considerations. For example, when considering the control of a plane, one can impose a constraint on the minimal altitude based on safety considerations. In this case, there is no actual interaction when the control system (plane) reaches the boundary of the state constraint. Moreover, the violations of the state constraint may be physically realizable and even admissible if it is sufficiently small. The analysis of control of this type for a dynamical system, when there is no interaction between the system and the boundary in the state space, yields a standard statement of an optimal control problem with state constraints [1]–[7], [10]– [12] which does not allow for the interaction with the boundary.

A typical optimal control problem of this type is the following *(P1)*:

$$\dot{x} = f(x,u), \quad u \in U;$$

$$x(0) = x_0, \quad x(T) = x_1;$$

$$x(t) \in G \quad for\ all \quad t \in [0,T], \quad G = \left\{x \in R^n : g(x) \le 0\right\};$$

$$J(x(\cdot),u(\cdot)) = \int_0^T f^0(x(t),u(t))\,dt \to \min.$$

Here it is assumed that $\dfrac{\partial g(x)}{\partial x} \ne 0 \quad \forall x : g(x) = 0$ and the set $f(x,U) = \bigcup_{u \in U} f(x,u)$ is convex.

The state constraints substantially complicate a problem. This is due to the fact that the state constraints introduce discontinuities into the dynamics of the system. Indeed it occurs because every time when considered trajectory of the control system reaches the boundary of the state constraints its set of all possible velocities must be intersected with the tangent cone to the set $G$ at this point and this operation involves the discontinuity in the dynamics of control system.

The situation is especially complicated when there is no any *a priory* information about the behavior of the optimal trajectory along the boundary of the state constraints. These difficulties are responsible for the appearance, as the Lagrange multipliers, measures and functions of bounded variations in necessary optimality conditions for problems with state constraints. The presence of

such general objects leads to qualitatively new phenomena that should be taken into account when deriving the corresponding necessary optimality conditions. In particular, the necessary optimality conditions with state constraints may degenerate when whey hold on any admissible trajectory. The new version of the Pontryagin maximum principle, which provides in some cases a possibility of complete investigation of this degeneracy phenomenon, was recently developed in [2]–[4].
Let

$$\mathrm{H}(x,u,\psi^0,\psi) = \left\langle f(x,u),\psi \right\rangle + \psi^0 f^0(x,u)$$

denote the Hamilton—Pontryagin function and

$$H(x,\psi^0,\psi) = \max_{u \in U} \mathrm{H}(x,u,\psi^0,\psi)$$

denote the Hamiltonian of the problem *(P)*.
The following result follows directly from [3].

*Theorem 1 (Maximum principle).* Let $x_*(t), u_*(t), T_*$ be an optimal triple in the problem *(P1)*. Then there exist a number $\psi^0 \geq 0$, an absolutely continuous function $\psi(t)$ and a scalar bounded regular nonnegative Borel measure $\eta$ on $[0, T_*]$ such that the following conditions hold:

a) Function $\psi(t)$ satisfy to the adjoint system

$$\dot{\psi}(t) \overset{a.e.}{=} -\frac{\partial \mathrm{H}}{\partial x}(x_*(t), u_*(t), \psi(t) + \int_{t_{1,0}}^{t} \frac{\partial g(x_*(t))}{\partial x} d\eta);$$

b) The maximum condition takes place:

$$\mathrm{H}(x_*(t), u_*(t), \psi(t) + \int_{t_{1,0}}^{t} \frac{\partial g(x_*(t))}{\partial x} d\eta) \overset{a.e.}{=} H(x_0(t), \psi(t) + \int_{t_{1,0}}^{t} \frac{\partial g(x_*(t))}{\partial x} d\eta);$$

c) For all $t \in [0, T_*]$ the following stationarity condition holds:

$$H(x_*(t), \psi(t) + \int_{0}^{t} \frac{\partial g(x_*(t))}{\partial x} d\eta) = H(x_*(t), \psi(t) + \int_{0}^{t} \frac{\partial g(x_*(t))}{\partial x} d\eta - \eta(t)) = 0;$$

Here $\eta(t)$ is an atomic component of the measure $\eta$ at the point $t$;

d) The nontriviality condition takes place

$$\left| \psi^0 \right| + \left\| \psi(0) \right\| + \left\| \eta \right\| \neq 0,$$

135

where

$$\|\eta\| = \sup_{\|x\|_{C[0,T_*]}=1} \int_0^{T_*} x(s)\,d\eta.$$

The main difference of the Theorem 1 from the other necessary optimality conditions for problems with state constraints [10], [11] consists in condition (c) of the stationarity of the Hamiltonian. This condition is quite natural and apparently should exist in one or another form of any "correctly" formulated necessary optimality conditions that generalize the Lagrange multipliers rule. The sense of this stationarity condition may be interpreted for some problems of mechanics as the energy conservation law.

Note that, the stationarity of the Hamiltonian in problems without state constraints follows from the other conditions of the maximum principle [12]. However, in the presence of state constraints, the derivation of this condition in the full form is associated with considerable difficulties (see [3], [6], [7] for details).

Consider now the following time-optimal problem with state constraints *(P2)*:

$$\dot{x} = f(x,u), \quad u \in U;$$

$$x(0) = x_0, \quad x(T) = x_1;$$

$$x(t) \in G \quad for\ all \quad t \in [0,T], \quad G = \{x \in R^n : g(x) \le 0\};$$

$$J(x(\cdot),u(\cdot)) = T \to \min.$$

Assume that there exist constants $\varepsilon_0 > 0, \alpha < 0$ and a locally Lipschitzian vector function $v : v(x) \in U(x), \forall x \in R^n$ such that

$$\left\langle \frac{\partial g(x)}{\partial x}, f(x,v(x)) \right\rangle \le \alpha \quad \forall x : 0 \le g(x) \le \varepsilon_0.$$

Under this controllability condition it is possible to approximate problem *(P2)* by a sequences of optimal control problems without state constraints.

For $i = 1,2,\ldots$ consider the following sequence of optimal control problems $(P_i)$:

$$\dot{x} = (1 - ih^2(x))f(x,u) + ih^2(x)f(x,v(x)), \quad u \in U;$$

$$x(0) = x_0, \quad x(T) = x_1;$$

$$J = T \to \min.$$

Here $h(x) = \max\{0, g(x)\}$. Obviously the function $h^2(x)$ is smooth and

$$\frac{\partial h^2(x)}{\partial x} = 2h(x)\frac{\partial g(x)}{\partial x} \quad \forall x \in R^n.$$

136

The constructed sequence of optimal control problems $(P_i)$, $i = 1,2,\ldots$ provides an approximation of the initial problem *(P2)*. The main properties of this approximation can be characterized by the following propositions.

*Proposition 1.* An arbitrary admissible trajectory $x(t)$ of the initial control system

$$\dot{x} = f(x,u), \quad u \in U$$

of the problem *(P)*, satisfying to the inequality $g(x(t)) \leq 0 \;\forall t \geq 0$, is simultaneously a trajectory of the modified control system

$$\dot{x} = (1 - ih^2(x))f(x,u) + ih^2(x)v(x), \quad u \in U$$

of the problem $(P_i)$ for any $i = 1,2,\ldots$. Moreover, if $i \geq \sqrt{\dfrac{1}{\varepsilon_0}}$ then any admissible trajectory $x(t)$ of the modified control system with initial condition $x_0 : g(x_0) \leq 0$ is also a trajectory of the initial control system and the following inequality holds: $ih^2(x(t)) < 1 \;\forall t \in I$.

*Proposition 2.* Let $\|f(x,U)\| = h(F(x),0) \leq M \;\forall x \in R^n$, where $M > 0$. Then for arbitrary trajectory $x(t)$ of the modified control system with initial condition $x_0 : g(x_0) \leq 0$ the following inequality holds:

$$\rho(x(t),G) \leq \frac{M}{\alpha}\sqrt{\frac{1}{i}} \quad \forall t \in [0,T].$$

*Proposition 3.* Let $\|F(x)\| \leq M \;\forall x \in R^n$, $M > 0$ and $T_*$ is an optimal time in a problem (P). Then for all $i = 1,2,\ldots$ there is a solution $x_i(t)$ of the problem $(P_i)$ with corresponding optimal time $T_i \leq T_*$. Further $T_i \to T_*$ as $i \to \infty$. The family of trajectories $\{x_i(t)\}$ is relatively compact in $C[0,T_*]$ and its arbitrary limit point $x(t)$ is an optimal trajectory in the problem (P).

## 3. The time-optimal problem of crossing a given domain.

Consider now the following optimal control problem *(P3)*:

$$\dot{x} = f(x,u), \quad u \in U;$$
$$x(0) = x_0, \quad x(T) = x_1;$$
$$J(x,T) = \int_0^T \delta(x)\, dt \to \min.$$

Here, as usually $x \in R^n$ is a state vector, $u \in R^m$ is a control parameter, $U$ is nonempty compact subset of $R^m$, $\delta$ is a characteristic function of the closed set $Z \subset \{x \in R^n : g(x) \le 0\}$, i.e.

$$\delta(x) = 1, \quad if \quad x \in Z;$$
$$\delta(x) = 0, \quad if \quad x \notin Z.$$

It is assumed that the vector function $f(x,u)$ and the scalar function $g(x)$ are smooth and $\dfrac{\partial g(x)}{\partial x} \ne 0$ if $g(x) = 0$. Further, there exists a constant $C_0 \ge 0$ such that

$$\langle x, f(x,u) \rangle \le C_0 (1 + \|x\|^2) \quad \forall x \in R^n, \forall u \in U;$$

and the set $F(x) = \bigcup \{f(x,u) : u \in U\}$ is convex for all $x \in R^n$.

Earlier some necessary optimality conditions for the problem *(P3)* were developed in [13], [14] under special conditions concerning the control system, possible behavior of an optimal trajectory and the fixed final time $T$.

In [13] the problem *(P3)* was considered in the case of linear control system and convex set $Z$. Moreover it was supposed in [13], that an optimal trajectory $x_*(t)$ has a final number of intersections with the boundary of the set $Z$. It was shown in [13] that in this case it is possible to reduce the original problem to the mathematical programming problem and by the use of the generalized Lagrange multipliers rule from [16] to develop the integral maximum principle. In [14] some special cases of the different collocation of the set $Z$ and points $x_0$ and $x_1$ which were not considered in [13]. In [15] the analogous results were obtained for a problem with delay.

The following recent result [8] is a necessary optimality conditions for the problem *(P3)* without any *a priori* assumptions concerning the optimal trajectory.

*Theorem 2 (Maximum Principle).* Let $u_*(t), x_*(t), T_*$ be an optimal triple. Then there exist a number $\psi^0 \le 0$, an absolutely continuous function $\psi(t)$ on $[0, T_*]$ and a bounded regular Borel measure $\eta$ on $[0, T_*]$: $\mathrm{supp}\, \eta \subset \{t \in [0, T_*] : g(x_*(t)) = 0\}$ such that the following conditions hold:

a) $\psi(t)$ is a solution to the adjoint system

$$\dot{\psi} \overset{a.e.}{=} -\left[ \frac{\partial f(x_*(t), u_*(t))}{\partial x} \right]^* \left( \psi - \int_0^t d\eta \right);$$

b) The maximum condition takes place:

$$H\left(x_*(t), \psi^0, \psi(t) - \int_0^t d\eta\right) \overset{a.e.}{=} \mathrm{H}\left(x_*(t), u_*(t), \psi^0, \psi(t) - \int_0^t d\eta\right);$$

c) The stationarity condition holds:

$$H(x_*(t), \psi_0, \psi(t) - \int_0^t d\eta) = H(x_*(t), \psi_0, \psi(t) - \int_0^t d\eta + \eta(t)) = 0 \quad \forall t \in [0, T_*],$$

here $\eta(t)$ is an atomic component of the measure $\eta$ at the point $t$;

d) There is a scalar bounded regular positive Borel measure $\nu$ on $[0, T_*]$ such that measure $\eta$ is absolutely continuous with respect to measure $\nu$ and $\nu$-a.e. the following equality holds:

$$\frac{d\eta}{d\nu} = \frac{\partial g(x_*(t))}{\partial x},$$

where $\dfrac{d\eta}{d\nu}$ is the Radon—Nikodim derivative;

e) The nontriviality condition holds

$$|\psi^0| + \|\psi(0)\| + \|\eta\| \neq 0.$$

## 4. Conclusion.

In the present paper two optimal control problems, which involve in their settings risky domains were considered. There are the optimal control problem with state constraints and the problem of time-optimal moving through the given domain. Both these problems involve essential discontinuities in respect to the state variables in the dynamic or in the functional. The important feature of these problems consists of the unknown and possibly complex character of moving of an optimal trajectory along the surface of problem discontinuity (the boundary of the risky domain $Z$). These possibly complex behaviors of an optimal trajectory produce the Lagrange multipliers of the very general nature (Borel measures) in the corresponding necessary optimality conditions of the first order.

In the cases of both problems due to their similar nature the developed necessary optimality conditions (theorems 1 and 2) are look also similar. It is important to underline that these results are developed by the use of the similar approximation technique. In both cases the initial nonclassical optimal control problems were approximated by sequences of the classical smooth optimal control problems with further limit procedure in the relations of Pontryagin maximum principle for approximating problems. So, these results are obtained as a consequence of the classical Pontryagin maximum principle [12] for smooth optimal control problems. Earlier the analogous approximation technique was used in [1]–[7].

The review of approximation methods of this type is given in [6].

## References

1. Arutyunov A.V., Perturbations of extremal problems with constraints, and necessary optimality conditions, Itogi Nauki i Tekhniki. Ser. Mat. Analiz. Vol. 27. Pp. 147–235. Moscow: VINITI. 1989. English transl. in J. of Soviet Math. Vol. 54. 1991.

2. Arutyunov A.V., Aseev S.M., The maximum principle for optimal control problems with state constraints. Nondegeneracy and stability, Dokl. Ross. Akad. Nauk. Vol. 334. Pp. 134–137. 1994. English transl. in Russian Acad. Sci. Dokl. Math. Vol. 49. 1994.

3. Arutyunov A.V., Aseev S.M., Investigation of the degeneracy phenomenon of the maximum principle for optimal control problems with state constraints, SIAM J. Control and Optimization. Vol. 35. No. 3. Pp. 930–952. 1997.

4. Arutyunov A.V., Aseev S.M., Blagodatskikh V.I., Necessary conditions of the first order in the problem of optimal control of a differential inclusion with phase constraints, Mat. Sb. Vol. 184. No. 6. Pp. 3–32. 1993. English transl. in Russian Acad. Sci. Sb. Math. Vol. 79. 1994.

5. Aseev S.M., A method of smooth approximations in the theory of necessary optimality conditions for differential inclusions, Izvestiya RAN: Ser. Mat. Vol. 61. No. 2. 1997. English transl. in Izvestiya: Mathematics, Vol. 61, No. 2, Pp. 235–258.

6. Aseev S.M., Methods of regularization in nonsmooth problems of dynamic optimization, J. of Mathematical Sciences. Vol. 94. No. 3. Pp. 1366–1393. 1999.

7. Aseev S.M., Extremal problems for differential inclusions with state constraints, Proceedings of the Steklov Institute of Mathematics. Vol. 233. Pp. 5–70. 2001.

8. Aseev S.M., Smirnov A.I., The necessary optimality conditions of the first order for the problem of optimal crossing a given domain. (In progress). 2001.

9. Clarke F.H., Optimization and nonsmooth analysis. New York: Wiley — Interscience. 1983.

10. Dubovitskii A.Ya., Milyutin A.A., Extremal problems with constraints, USSR Comput. Math. Phys. 1965, (5), 1–80.

11. Gamkrelidze R.V., Optimal control processes with bounded phase coordinates, USSR Math. Izvestiya, 1960, (24), 315–356.

12. Pontryagin L.S., Boltyanskii B.G., Gamkrelidze R.V., and E.F. Mischenko, The mathmatical theory of optimal processes. New York: Interscience Publishers, John Wiley & Sons, Inc. 1962.

13. Pshenichnyi B.N., Ochilov S., On a problem of optimal moving through a given domain, Cybernetics and Computer Technique, 1993, (99), 3–8.

14. Pshenichnyi B.N., Ochilov S., On a special time-optimal problem, Cybernetics and Computer Technique, 1994, (101), 11–15.

15. Otakulov S., Ochilov S., On time-optimal problem of moving through the given domain, in Abstracts of the reports of the Voronezh spring mathematical school "Pontryagin's studies — ", Contemprorary methods in the theory of the boundary problems, Voronezh, May, 3–9, 1999. Voronezh, 1999.

16. Pshenichnyi B.N., Necessary conditions of extremum. Moscow: Nauka. 1982.