

Geostatistical predictions based on block data

Shigeru Mase¹, Kouji Takahashi, and Nobuyuki Mochida

Tokyo Institute of Technology
Meguro-ku, Tokyo, Japan
(e-mail: mase@is.titech.ac.jp)

Abstract. Block-to-block and block-to-point kriging predictions based on block data are proposed. Blocks may be regular (mesh data) or of more general shapes. Under the assumptions of second-order stationarity and isotropicity, we show how to lessen the number of calculations of relevant block-to-block covariances. As illustrations, a mesh data of population and a simulated block data on convex polygons are analyzed.

Keywords: Geostatistics, Kriging predictor, Block data, Block-to-point prediction, Mesh data, Change of scale problem, Gaussian Random Field.

1 Introduction

Geostatistics has the origin in the pioneering work of South African mining engineer D. G. Krige in 1950's who introduced a statistical methodology to evaluate gold ore grade based on boring core samples. In 1970's, French mathematician G. Matheron formulated a regression based spatial prediction method for which he coined the term "kriging". Although it has been developed mainly in application fields and outside of the usual statistical community, now kriging method has become an indispensable statistical tool in variety of fields such as epidemiology, environment science, ecology, agriculture, geology, civil engineering, social sciences, geography, fishery science, oceanography and so on where available data are only small portions of a large spatial structure and one want to know the global spatial distribution of a feature.

The probabilistic basis of geostatistics is a random field $Z(\mathbf{x})$, $\mathbf{x} \in \mathbf{R}^d$, d being typically two or three. Available data is a set of observations $Z(\mathbf{x})$ at specified locations $\mathbf{x} = \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ and one wants to predict the value $Z(\mathbf{x}_0)$ for arbitrary locations \mathbf{x}_0 . Hence it can be thought as a spatial interpolation and/or extrapolation of data. By the way, the term "estimation" instead of "prediction" has been frequently used in geostatistics literature from a historical reason.

As well as predicting a point value $Z(\mathbf{x}_0)$ (point kriging), it is sometimes required to predict a block value $Z(B) = |B|^{-1} \int_B Z(\mathbf{x}) d\mathbf{x}$ (block kriging) which is the mean of $Z(\mathbf{x})$ over a block B . As to these kriging problems, there are well-established results, see, e.g. Chiles and Delfiner[2], Cressie[3] and Wackernagel[5]. In this paper, we will discuss the converse problem, that is, kriging predictions of point or block values based on block data. The use of

data of type $Z(B)$ rather than original $Z(\mathbf{x})$ is sometimes called “the change of support problem” in the literature and known to cause various problems, see Cressie[4] and Chiles and Delfiner[2].

2 Second-order stationary random fields

In the following, random fields $Z(\mathbf{x})$ are assumed to be second-order stationary, that is, the mean $\mathbf{E}\{Z(\mathbf{x})\}$ is a constant μ irrespective of \mathbf{x} and the covariance $\text{Cov}\{Z(\mathbf{x}), Z(\mathbf{y})\}$ is a function $C(\mathbf{x} - \mathbf{y})$ of the difference $\mathbf{x} - \mathbf{y}$ only. C is called the covariance function of Z . It is even and is characterized by the positive-definiteness

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j C(\mathbf{x}_i - \mathbf{x}_j) \geq 0$$

for any $\{\mathbf{x}_i\}$ and constants $\{c_i\}$, where the equality holds only for $c_1 = c_2 = \dots = c_n = 0$. A second-order stationary random field is said to be isotropic if its covariance function is the function of the norm $|\mathbf{x} - \mathbf{y}|$ of the difference $\mathbf{x} - \mathbf{y}$.

The followings are three typical isotropic covariance functions. They have two positive parameters a, b . The exponential covariance family is $C_{exp}(\mathbf{h}) = b \exp(-|\mathbf{h}|/a)$.

$$C_{sph}(\mathbf{h}) = \begin{cases} b(1 - 3|\mathbf{h}|/(2a) + |\mathbf{h}|^3/(2a^3)), & |\mathbf{h}| \leq a, \\ 0, & |\mathbf{h}| > a. \end{cases}$$

And the Gaussian covariance family is $C_{gau}(\mathbf{h}) = b \exp(-|\mathbf{h}|^2/a)$.

In geostatistics, the concept of intrinsic stationarity has been preferred to second-order stationarity. A random field is intrinsic stationary if

$$\mathbf{E}\{Z(\mathbf{x}) - Z(\mathbf{y})\} = 0, \quad \mathbf{E}\{|Z(\mathbf{x}) - Z(\mathbf{y})|^2\} = 2\gamma(\mathbf{x} - \mathbf{y})$$

for all \mathbf{x}, \mathbf{y} . The even function γ which depends only on the difference $\mathbf{x} - \mathbf{y}$ is called the (semi)variogram which may be unbounded contrary to covariance functions. The use of variograms is intended to cancel a possible linear trend which seems frequent in mining data. The concept of intrinsic stationarity is more general than second-order stationarity in principle since it does not assume the existence of the mean and variance of $Z(\mathbf{x})$. If Z is second-order stationary, it is intrinsic stationary and the relation $\gamma(\mathbf{h}) = C(\mathbf{0}) - C(\mathbf{h})$ holds. In variogram based kriging, it is usual to estimate model parameters by fitting theoretical variograms to binned sample variograms using the method of least squares. In order to apply this method, it is essential to know distances between locations of each pair of data. But it is impossible or ambiguous to define distances between two blocks. This is the main reason why we use the second-order stationarity assumption and the

maximum likelihood approach to estimate model parameters assuming the normality of random fields in this paper. This has further advantages that we need not classify data and can reduce three parameters models to one parameter models, see Prop. 3.

3 Ordinary kriging for block data

Let $Z = Z(\mathbf{x})$ be a second-order stationary random field with a covariance function $C(h)$ and a mean μ . Let B be a block (i.e., a bounded region with positive volume). The block data of Z for the block B is defined as follows

$$Z(B) = \frac{1}{|B|} \int_B Z(\mathbf{x}) d\mathbf{x},$$

where $|B|$ is the volume of B . The mean of $Z(B)$ is μ . Let B_1, B_2, \dots, B_n be a set of blocks. They are not necessary disjoint. The covariance between two block data $Z(B_\alpha)$ and $Z(B_\beta)$ is given by

$$\begin{aligned} C_{B_\alpha, B_\beta} &= \frac{1}{|B_\alpha||B_\beta|} \iint_{B_\alpha \times B_\beta} \text{Cov}\{Z(\mathbf{x}), Z(\mathbf{y})\} d\mathbf{x} d\mathbf{y} \\ &= \frac{1}{|B_\alpha||B_\beta|} \iint_{B_\alpha \times B_\beta} C(\mathbf{x} - \mathbf{y}) d\mathbf{x} d\mathbf{y}. \end{aligned}$$

Let B_0 be a new block and we want to predict $Z(B_0)$ by a linear combination of block data $Z_B = (Z(B_1), Z(B_2), \dots, Z(B_n))^T$:

$$\hat{Z}(B_0) = \sum_{\alpha=1}^n w_\alpha Z(B_\alpha).$$

This block-to-block kriging prediction can be constructed according to the standard procedure of the ordinary kriging based on point data as explained in Cressie[3], Wackernagel[5], or Chiles and Delfiner [2]. In order to guarantee the unbiasedness, we put the constraint $\sum_i w_i = 1$ on weights w_1, w_2, \dots, w_n (ordinary kriging). Hence

$$\mathbf{E} \left\{ \hat{Z}(B_0) \right\} = \sum_{\alpha=1}^n w_\alpha \mathbf{E} \{ Z(B_\alpha) \} = \mu \sum_{\alpha=1}^n w_\alpha = \mu.$$

Its mean squared prediction error $\sigma_E^2 = \mathbf{E} \left\{ (\hat{Z}(B_0) - Z(B_0))^2 \right\}$ is

$$\sum_{\alpha=1}^n \sum_{\beta=1}^n w_\alpha w_\beta C_{B_\alpha, B_\beta} + C_{B_0, B_0} - 2 \sum_{\alpha=1}^n w_\alpha C_{B_\alpha, B_0}.$$

Since we have to minimize σ_E^2 under the constraint $\sum_{\alpha} w_{\alpha} = 1$, consider the objective function with the Lagrange multiplier λ :

$$\phi(\{w_{\alpha}\}, \lambda) = \sigma_E^2 - 2\lambda \left(\sum_{\alpha=1}^n w_{\alpha} - 1 \right).$$

Proposition 1. *Weights $\{w_{\alpha}\}$ of the ordinary block-to-block kriging predictor $\hat{Z}(B_0)$ are the solution of the following equation:*

$$\begin{pmatrix} C_{B_1, B_1} & \cdots & C_{B_1, B_n} & 1 \\ \vdots & \ddots & \vdots & \vdots \\ C_{B_n, B_1} & \cdots & C_{B_n, B_n} & 1 \\ 1 & \cdots & 1 & 0 \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_n \\ -\lambda \end{pmatrix} = \begin{pmatrix} C_{B_1, B_0} \\ \vdots \\ C_{B_n, B_0} \\ 1 \end{pmatrix}$$

where λ is the Lagrange multiplier. This system is authentic, i.e., $\hat{Z}(B_i) = Z(B_i)$ for $i = 1, 2, \dots, n$. The corresponding mean squared prediction error is $\sigma_E^2 = \lambda + C_{B_0, B_0} - \sum_{\alpha=1}^n w_{\alpha} C_{B_{\alpha}, B_0}$.

In a similar way, we can consider the block-to-point ordinary kriging. The covariance between block and point data $Z(B)$ and $Z(\mathbf{y})$ is given by

$$C_{B, \mathbf{y}} = \frac{1}{|B|} \int_B \text{Cov}\{Z(\mathbf{x}), Z(\mathbf{y})\} d\mathbf{x} = \frac{1}{|B|} \int_B C(\mathbf{x} - \mathbf{y}) d\mathbf{x}.$$

The block-to-point ordinary kriging predictor of $Z(\mathbf{x}_0)$ takes the form $\hat{Z}(\mathbf{x}_0) = \sum_{\alpha=1}^n w_{\alpha} Z(B_{\alpha})$ with constraint $\sum_i w_i = 1$.

Proposition 2. *Weights $\{w_{\alpha}\}$ of the ordinary block-to-point kriging predictor $\hat{Z}(\mathbf{x}_0)$ are the solution of the following equation:*

$$\begin{pmatrix} C_{B_1, B_1} & \cdots & C_{B_1, B_n} & 1 \\ \vdots & \ddots & \vdots & \vdots \\ C_{B_n, B_1} & \cdots & C_{B_n, B_n} & 1 \\ 1 & \cdots & 1 & 0 \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_n \\ -\lambda \end{pmatrix} = \begin{pmatrix} C_{B_1, \mathbf{x}_0} \\ \vdots \\ C_{B_n, \mathbf{x}_0} \\ 1 \end{pmatrix}$$

where λ is the Lagrange multiplier. The corresponding mean squared prediction error is $\sigma_E^2 = \lambda + C(\mathbf{0}) - \sum_{\alpha=1}^n w_{\alpha} C_{B_{\alpha}, \mathbf{x}_0}$.

4 Estimation of model parameters

As to kriging model parameter estimations, there are two main methods. One is the traditional least square fitting of theoretical variograms to sample variograms and the other is the maximum likelihood estimation assuming the normality of data. In order to use least square fittings, one needs to define distances between blocks, but this is impossible or ambiguous for blocks.

So we employ the maximum likelihood estimation assuming Z is a Gaussian random field. A merit of this approach is one can use original data themselves directly. On the other hand, variogram based fittings have to change data into suitable class means.

Proposition 3. *Let Z be a stationary Gaussian random field with mean μ and covariance function $bC(a^{-1}\mathbf{x})$, $a, b > 0$. Let $Z_B = (Z(B_1), Z(B_2), \dots, Z(B_n))$ be a block data vector. Then the maximum likelihood estimators \hat{a} , \hat{b} and $\hat{\mu}$ satisfy the relations*

$$\mu = \frac{\mathbf{e}^T \Sigma_B^{-1}(a, 1) Z_B}{\mathbf{e}^T \Sigma_B^{-1}(a, 1) \mathbf{e}}, \tag{1}$$

$$b = \frac{1}{n} (Z_B - \mu \mathbf{e})^T \Sigma_B^{-1}(a, 1) (Z_B - \mu \mathbf{e}), \tag{2}$$

$$0 = \{ \Sigma_B^{-1}(a, 1) (Z_B - \mu \mathbf{e}) \}^T \frac{d \Sigma_B(a, 1)}{da} \{ \Sigma_B^{-1}(a, 1) (Z_B - \mu \mathbf{e}) \} - b \operatorname{tr} \left(\Sigma_B^{-1}(a, 1) \frac{d \Sigma_B(a, 1)}{da} \right) \tag{3}$$

where $\mathbf{e} = (1, 1, \dots, 1)^T$ and $\Sigma_B(a, b)$ is the covariance matrix of Z_B . In particular, (3), after μ and b being eliminated using (1) and (2), is an equation of variable a only and we have \hat{a} by solving this equation. Then $\hat{\mu}$ and \hat{b} can be calculated immediately from relations (1) and (2).

5 Mesh data case

In order to apply block data kriging and model parameter estimations, it is essential to compute block covariance matrix efficiently. In general, this is difficult and time-consuming if not impossible. If there is n blocks, we need to calculate $n(n + 1)/2$ covariances numerically in principle.

Many spatial data are given as aggregates of original data per mesh. Typical examples are demographic data. For mesh-type blocks, we can reduce the number of necessary computations using stationarity and isotropy.

Proposition 4. *If B_1 and B_2 are disjoint, $C_{A, B_1 \cup B_2} = C_{A, B_1} + C_{A, B_2}$. Let Z be a stationary random fields. Then $C_{A, B} = C_{A+\mathbf{h}, B+\mathbf{h}}$. If, moreover, Z is isotropic and T is a congruent transformation, $C_{A, B} = C_{T(A), T(B)}$.*

If B_1, B_2, \dots, B_n , $n = n_1 \times n_2$, are rectangles which consists of n_1 by n_2 congruent division of a rectangle, we need to compute only $n_1 + n_2 - 1 + 2(n_1 - 1)(n_2 - 1)$ covariances C_{B_i, B_j} instead of $n(n + 1)/2$ ones if Z is second-order stationary. If, moreover, Z is isotropic, we need to compute only $n_1 + n_2 - 1 + (n_1 - 1)(n_2 - 1)$ covariances. For example, if $n_1 = n_2 = 10$, we need to compute only 100 and 181 covariances respectively instead of 5,500 ones.

Proposition 5. *If Z is a two-dimensional second-order stationary random field and A is a rectangle with width s and height t ,*

$$C_{A,A+\mathbf{h}} = |A|^{-2} \int_{-s}^s \int_{-t}^t (s - |x|)(t - |y|)C((x, y)^T - \mathbf{h})dxdy.$$

As an application, we apply the block-to-point kriging to a 10 by 10 mesh data, which are populations (1,000/ km^2) of Tokyo metropolitan area (Japanese Statistics Bureau (2000)). Strictly speaking, this is not a block data but count data per mesh. We assume a hypothetical field of population density and pretend this is a resulting block data. We fitted spherical, exponential and Gaussian models. Fig. 1 shows resulting contour images of block-to-point kriging predictions. Three results show fairly similar features.

6 General regions case

Computation of block-to-block or block-to-point covariance matrices is quite difficult if block shapes are arbitrary. In this section, we propose an algorithm of computing these covariances approximately. We assume that Z is second-order stationary and isotropic.

Proposition 6. *Assume Z is a two-dimensional stationary and isotropic random field. Then the block-to-block covariance $C_{A,B}$ is*

$$C_{A,B} = \int_{r_0}^{r_1} G_{A,B}(r)C(r)rdr,$$

where

$$G_{A,B}(r) = \frac{1}{|A||B|} \int_0^{2\pi} |(A - re^{i\theta}) \cap B|d\theta. \quad (4)$$

$r_0 \geq 0$ (resp. r_1) is the minimum (resp. maximum) of the set $\{r \geq 0 : (A - re^{i\theta}) \cap B \neq \emptyset \exists \theta\}$. r_1 is always finite.

Also the block-to-point covariance $C_{A,\mathbf{y}}$ is $C_{A,\mathbf{y}} = \int_{r_0}^{r_1} G_{A,\mathbf{y}}(r)C(r)rdr$ where $G_{A,\mathbf{y}}(r) = \frac{1}{|A|} \int_0^{2\pi} \mathbf{1}_A(\mathbf{y} + re^{i\theta})d\theta$ and $r_0 \geq 0$ (resp. r_1) is the minimum (resp. maximum) of the set $\{r \geq 0 : \mathbf{y} + re^{i\theta} \in A \exists \theta\}$. r_1 is again finite.

As an application, we show a simulation result for block data kriging for convex regions. The basic region D is the square $[0, 10] \times [0, 10]$ and we generated Voronoi cells B_i with centers generated using a simple sequential inhibition point process. Actually those 83 Voronoi cells completely included in D were used, see Fig. 2. Gaussian random fields were generated on D with the exponential model with parameters $b = 10$, $a = 1, 2, 3$ and the spherical model with parameters $b = 10$, $a = 4, 5$ and block data approximated by discrete sums were generated.

Both $G_{A,B}(r)$ and $G_{A,\mathbf{y}}(r)$ have no simple closed expressions in general. A practical procedure is to compute their values at sufficiently many r 's and interpolate them. Since these functions depend only on A, B (resp. A, \mathbf{y}) and does not depend on covariance functions $C(r)$, we need compute them only once. Also the area $|(A - re^{i\theta}) \cap B|$ in (4) needs two-dimensional integrations over irregular sets which can be efficiently evaluated by a quasi-Monte Carlo integration using low-discrepancy sequences.

Fig. 2 shows used Voronoi cells (top left), the original random field image (top right), the corresponding block image (bottom left), and the block-to-point kriging prediction result (bottom right).

7 Conclusion

Block data kriging, in particular, block-to-point kriging seems useful since many data such as in demography and epidemiology are often publicized as aggregates per municipalities or as mesh data from the first. Such data may be also analyzed using so-called hierarchical Bayes models with Markov Chain Monte Carlos as explained in detail in [1]. Applicability of this method may be more general than the present one since it does not necessary assume a stationary random field framework. On the other hand, it requires a data specific hierarchical Bayes model.

It should be borne in mind that block data are smoothed from the first and, therefore, one cannot expect to recover finer details of the original data. Also blocks with too irregular shapes may lessen the discriminative power of covariance matrices apart from numerical inefficiency.

In order to apply block-to-block and block-to-point kriging, one has to compute a lot of block covariances efficiently. In this paper, we showed that this is feasible at least for the two-dimensional second-order isotropic and isotropic case. For non-isotropic cases, parallel computations are probably the last resort.

References

1. Banerjee, S., Carlin, B. P. & Gelfand, A. E. (2003). *Hierarchical Modeling and Analysis for Spatial Data*, London: Chapman & Hall/CRC.
2. Chiles, J.-P. & Delfiner, P. (1999). *Geostatistics: Modeling Spatial Uncertainty*, New York: Wiley & Sons.
3. Cressie, N. A. C. (1991). *Statistics for Spatial Data*, New York: Wiley & Sons.
4. Cressie, N. A. C. (1996). Change of support and the modifiable areal unit problem, *Geographical Systems*, 3, 159-180.
5. Wackernagel, H. (1999). *Multivariate Geostatistics: An Introduction With Applications* (2nd. Rev.), New York: Springer Verlag.

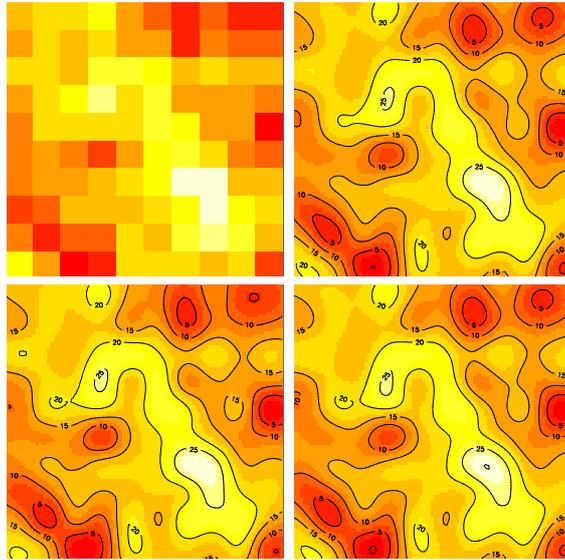


Fig. 1. Mesh data of populations ($1,000/km^2$) (top left). Block-to-point kriging result using the spherical model (top right), the exponential model (bottom left), and the Gaussian model (bottom right).

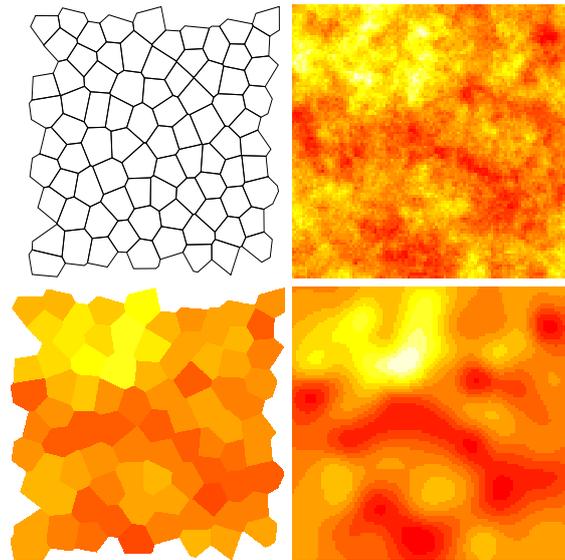


Fig. 2. Voronoi cells (top left), the original random field image (top right), the block data (bottom left) and the block-to-point kriging result (bottom right). Note that color levels for three images are slightly different.

Drag-Free, Attitude and Formation Control for Monitoring the Earth Gravity Field

Luca Massotti¹, Andres Molano² and Enrico Canuto²

¹ESA ESTEC, Earth Observation - Future Missions Division (EOP-SFP),
2200 AG, Noordwijk ZH, The Netherlands

E-mail: luca.massotti@esa.int

²Politecnico di Torino, Dipartimento di Automatica e Informatica, Corso Duca degli
Abruzzi 24, 10129 Torino, Italy

E-mail: enrico.canuto@polito.it

Abstract: Following the successful launch in Spring 2009 of the GOCE satellite (Gravity field and steady-state Ocean Circulation Explorer), a more ambitious mission consisting in a satellite formation of two satellite separated by (a minimum) 10 km distance, is under study at the European Space Agency, aiming at monitoring the Earth's gravity field fluctuations, during a (at least) 6-year mission. Since GOCE is the first flying drag-free satellite, the envisaged formation might be the first drag-free formation, posing a suite of challenging technology and control problems under study and solution. The paper concentrates on a triad of control problems to be solved and coordinated (formation, drag-free and attitude), all of them being constrained by a long-life low-Earth-orbit mission, imposing low propellant mass, scarce electric propulsion throttability and limited electric power. Driving requirements are presented and discussed, showing how they can be met through Embedded Model Control design. Finally, realistic simulated results are included.

Keywords: Satellite formation, control, drag-free, attitude, low-Earth-orbit, gravity monitoring,

1 Introduction

One of the possible future Earth gravity monitoring missions after GOCE (Gravity field and steady-state Ocean Circulation Explorer), recently launched and successfully operating (Canuto, 2008, Canuto and Massotti, 2009, Canuto, Massotti and Molano, 2010), will be based on laser interferometry, in order to extend the gravity-gradient baseline to tens of km. A formation of at least two satellites is needed to implement long-baseline interferometry. In addition, a long mission is desirable to complement gravity spatial variations with time, and the orbit must be sufficiently low-altitude, to reveal high-order gravity harmonics. A mission of this kind is under study at the European Space Agency: in the last study (at the moment of this publication) a 10-km baseline and a mission length of 6 years have been selected. To allow scientific advancements, each satellite shall be drag-free, implying the residual CoM (Centre-of-Mass) non gravitational acceleration to be lower than $0.01 \mu\text{m/s}^2$ in a frequency band from 1 to 10 mHz. Similarly, proportionate requirements apply to angular accelerations, angular rates and attitude. Residuals are progressively relaxed below and above the mid-frequency band. In addition, a 3D formation must be kept, with loose requirements

at first glance: variations of the relative formation position must stay in a box $500 \times 50 \times 50 \text{ m}^3$ wide, the sequence of coordinates being along-track, cross-track and radial. Several technology problems have to be solved including propulsion, since the latter, even if essential for formation and drag-free, must be employed for attitude control too. Electric propulsion is mandatory in order to reduce propellant mass around 10% of the satellite mass (500 kg). Second, throttability (the max/min thrust ratio) must be sufficiently high to cope with a highly variable drag imposed by long-term and short-term solar activity. Already-flown, scalable though with insufficient throttability, micro-RIT (radio-frequency) thruster technology (Loeb, Schartner, Weiss, Feili and Meyer, 2004) is under study and test at Thales Alenia Space Italia premises. Since throttability looks one of the most critical technology constraints, control strategies must be designed so as to minimize thrust peak. Thruster layout, sketched in Fig. 1, must repeat the early GOCE design (Canuto and Massotti, 2009).

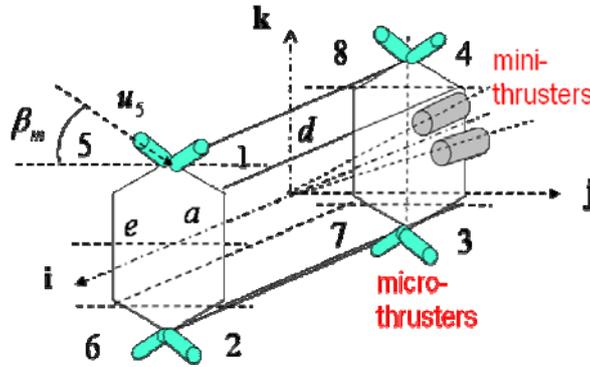


Fig. 1. Satellite shape and thruster layout.

A pair of larger thrusters (mini-thrusters, 0.4 mN to 18 mN range), in cold redundancy like on-board GOCE, will take care of along track drag-free and formation control (in \hat{i} direction, see Fig. 1), while eight smaller thrusters (micro-thrusters, 0.05 mN to 2 mN range) will accomplish lateral formation and drag-free control (\hat{j} and \hat{k} directions in Fig. 1), as well as attitude control for a total of 5 degrees-of-freedom (DoF).

Four classes of sensors are available on-board.

1. A pair of GPS receivers (1-Hz data rate) on each satellite will be employed for formation control and attitude reference generation. They should also be employed to calibrate accelerometer bias, as mentioned below.
2. A pair of GOCE-type accelerometers (10-Hz data rate) will be the sensors of the CoM drag-free control, and the wide-band attitude sensors for the attitude control, on each satellite. Accelerometer bias and drift are the main sources of formation separation (about 500 m at day), to be counteracted by formation control. Since the bias corresponds to $0.06 \div 6 \text{ mN}$ thrust range, the larger

value looks incompatible with micro-thruster range, thus asking either re-design or on-board calibration.

3. Two star trackers (2-Hz data rate) in cold redundancy for aligning attitude to the orbital reference frame. Since a single star tracker does not guarantee 3D uniform error, some problems can arise in attitude control.
4. On-board optical metrology, equally replicated on each satellite. Optical metrology allows to accurately measure distance variations along the optical interferometer baseline, as well as the 2D tilt of the satellite along-track axis (\vec{i}) with respect to the optical baseline. Consequently, the lateral displacement may be real-time monitored with the help of the attitude obtained from star trackers, and the optical distance from the metrology itself. In this way, an alternative metrology with respect to the differential GPS will be capable of providing formation relative position.

The paper is devoted to give an overview of the control strategies, namely 3D formation, drag-free and attitude. First, reference frames and satellite sensor and actuator dynamics are briefly reported, paying attention to disturbance and measurement error classes. Then control requirements and design are outlined. The paper ends with the most significant simulated results.

Control strategies are designed within the Embedded Model Control framework (Canuto, 2007), where control algorithms are built around a real-time Embedded Model, and split into reference generator, noise estimator and control law. Key to noise estimator is the definition of the noise channels (Canuto, Massotti and Molano, 2010). Noise estimator and Embedded Model may be interpreted as state observers. Embedded Model and control design are directly tackled in the discrete-time domain. Here continuous time is adopted.

2 Frames and dynamics

Dynamics is provided in the simplified form suitable to Embedded Model. The mean Local Orbital Reference Frame (LORF) $\mathcal{R}_o = \{C, \vec{i}_o, \vec{j}_o, \vec{k}_o\}$, centred in the formation CoM C , is defined by the instantaneous orbit orientation $\vec{v}/|\vec{v}|$, \vec{v} being the CoM velocity, and by the orbital plane orthogonal to the normalized angular momentum $\vec{h} = \vec{r} \times \vec{v}$, where $\vec{r} = (\vec{r}_0 + \vec{r}_1)/2$ is the formation CoM under equal satellite masses. Each satellite CoM position is denoted with \vec{r}_k , where $k = 0$ refers to the leader and $k = 1$ to the follower.

The LORF is the reference frame for science and attitude control: LORF axes are defined by

$$\vec{i}_o = \vec{v}/|\vec{v}|, \vec{j}_o = \vec{r} \times \vec{v}/|\vec{r} \times \vec{v}|, \vec{k}_o = \vec{i}_o \times \vec{j}_o. \quad (1)$$

Axes from \vec{i}_o to \vec{k}_o are respectively referred as along-track, cross-track and radial. Dropping arrows when inertial coordinates are considered, the matrix $R_o = [\vec{i}_o \quad \vec{j}_o \quad \vec{k}_o]$, directly obtained from GPS measurements, accomplishes the LORF-to-inertial coordinate transformation, and defines a common reference attitude to be tracked by both spacecrafts ($k = 0, 1$) during all over the mission.

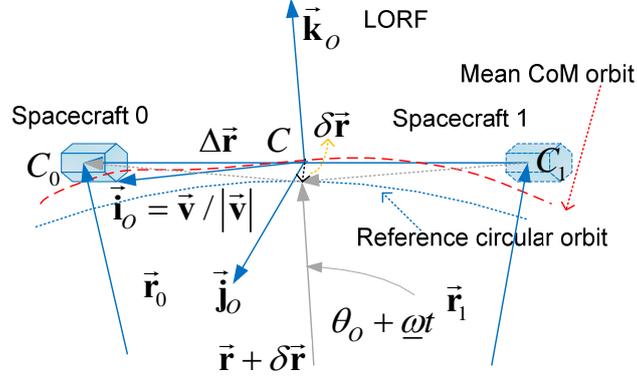


Fig. 2. Formation geometry and local orbital frame.

LORF dynamics may be either described through CoM dynamics, or through LORF quaternion \mathbf{q}_o , the latter directly obtained from R_o . Treating orbit angular rate $\boldsymbol{\omega}_o$ and acceleration $\mathbf{a}_o = \dot{\boldsymbol{\omega}}_o$ as state variables, the following LORF equations hold

$$\begin{aligned} \dot{\mathbf{q}}_o(t) &= \frac{1}{2} \mathbf{q}_o(t) \otimes \boldsymbol{\omega}_o(t), \mathbf{q}_o(0) = \mathbf{q}_{o0} \\ \dot{\boldsymbol{\omega}}_o(t) &= \mathbf{a}_o(t) + \mathbf{w}_o(t), \boldsymbol{\omega}_o(0) = \boldsymbol{\omega}_{o0}, \\ \dot{\mathbf{a}}_o(t) &= \mathbf{s}_o(t), \mathbf{a}_o(0) = \mathbf{a}_{o0} \\ \mathbf{y}_o(t) &= \mathbf{q}_o(t) \otimes \mathbf{e}_o(t) \end{aligned} \quad (2)$$

where \otimes denotes quaternion product, \mathbf{w}_o is a wide-band noise (white noise in discrete-time domain), \mathbf{s}_o is the angular jerk, and \mathbf{y}_o is the orbit quaternion measurement obtained through (1) from GPS range and range rate less the error quaternion \mathbf{e}_o . Feedback of the error quaternion to \mathbf{w}_o and \mathbf{s}_o allows to recover quaternion estimate together with orbit angular rate and acceleration, thus providing reference attitude trajectories. An equation similar to (2) applies to body quaternion \mathbf{q}_k ($k=0,1$), retrieved from the inertial coordinates of the body axes collected in $R_k = [\mathbf{i}_k \ \mathbf{j}_k \ \mathbf{k}_k]$ (see Fig. 1):

$$\begin{aligned} \dot{\mathbf{q}}_k(t) &= \frac{1}{2} \mathbf{q}_k(t) \otimes \boldsymbol{\omega}_k(t), \mathbf{q}_k(0) = \mathbf{q}_{k0} \\ \dot{\boldsymbol{\omega}}_k(t) &= J_k^{-1} (-\boldsymbol{\omega}_k(t) \times J_k \boldsymbol{\omega}_k(t)) + \mathbf{u}_{qk}(t) + \mathbf{a}_{qk}(t) + \mathbf{w}_{qk}(t), \boldsymbol{\omega}_k(0) = \boldsymbol{\omega}_{k0}, \\ \dot{\mathbf{a}}_{qk}(t) &= \mathbf{s}_q(t), \mathbf{a}_{qk}(0) = \mathbf{a}_{qk0} \\ \mathbf{y}_{qk}(t) &= \mathbf{q}_k(t - \tau_s) \otimes \mathbf{e}_{qk}(t - \tau_{sk}) \end{aligned} \quad (3)$$

where $\boldsymbol{\omega}_k$, \mathbf{a}_{kq} , \mathbf{s}_{kq} and \mathbf{w}_q have the same meaning as in Eq. (2), with the exception that now \mathbf{a}_{kq} accounts for un-modelled angular accelerations (Canuto, 2008). Command torques provided by thrusters are into the command acceleration vector \mathbf{u}_{kq} . J_k is the inertia matrix, close to be diagonal, but largely unbalanced because of a slender spacecraft as GOCE. Quaternion is retrieved from the star

tracker quaternion \underline{q}_{qk} less the error \underline{e}_{qk} and a delay τ_{sk} . The same equation applies to each satellite upon different notations. Attitude control is actuated ideally

$$\underline{q}_k = \underline{q} = \underline{q}_O, \quad (4)$$

less a tracking error

$$\underline{e}_k = \underline{q}_k^* \otimes \underline{q} = \begin{bmatrix} \underline{e}_{0k} \\ \underline{e}_k \end{bmatrix}. \quad (5)$$

Spacecraft and formation CoM dynamics may be written using a manipulated version of Hill's equation (Canuto, Massotti and Molano, 2010, Inalhan, Tillerson and How, 2002). With reference to Fig. 2, let us denote the Cartesian coordinates of the spacecraft k in the LORF frame with $\Delta \mathbf{r}_k$ and the (local) rate with $\Delta \mathbf{v}_k$. Due to LORF rotation, the following relative dynamics applies

$$\begin{aligned} \Delta \dot{\mathbf{r}}_k(t) &= \Delta \mathbf{v}_k(t), \quad \Delta \mathbf{r}_k(0) = \Delta \mathbf{r}_{k0} \\ \Delta \dot{\mathbf{v}}_k(t) &= -\dot{\boldsymbol{\omega}}_O \times \Delta \mathbf{r}_k - \boldsymbol{\omega}_O \times (\boldsymbol{\omega}_O \times \Delta \mathbf{r}_k + 2\Delta \mathbf{v}_k) \\ &\quad - \nabla \mathbf{g}(\mathbf{r}) \Delta \mathbf{r}_k + R_k(\underline{e}_k)(\mathbf{u}_k + \mathbf{d}_k + \mathbf{w}_k), \quad \Delta \mathbf{v}_k(0) = \Delta \mathbf{v}_{k0}, \\ \dot{\mathbf{d}}_k(t) &= \mathbf{s}_k(t), \quad \mathbf{d}_k(0) = \mathbf{d}_{k0} \end{aligned} \quad (6)$$

where gravity acceleration is reduces to the tidal component $\nabla \mathbf{g}(\mathbf{r}) \Delta \mathbf{r}_k$, and non gravitational accelerations have been split into command \mathbf{u}_k , disturbance \mathbf{d}_k and noise \mathbf{w}_k (Canuto, 2008). Transformation $R_k(\underline{e}_k)$ maps body coordinates into LORF. Now defining the LORF formation coordinate as

$$\Delta \mathbf{r} = \Delta \mathbf{r}_0 - \Delta \mathbf{r}_1, \quad (7)$$

and likely the differential rate $\Delta \mathbf{v}$, non gravitational acceleration $\Delta \mathbf{d}$, command $\Delta \mathbf{u}$ and noise $\Delta \mathbf{w}$, the formation equation may be written

$$\begin{aligned} \Delta \dot{\mathbf{r}}(t) &= \Delta \mathbf{v}(t), \quad \Delta \mathbf{r}(0) = \Delta \mathbf{r}_0 \\ \Delta \dot{\mathbf{v}}(t) &= -\dot{\boldsymbol{\omega}}_O \times \Delta \mathbf{r} - \boldsymbol{\omega}_O \times (\boldsymbol{\omega}_O \times \Delta \mathbf{r} + 2\Delta \mathbf{v}) \\ &\quad - \nabla \mathbf{g}(\mathbf{r}) \Delta \mathbf{r} + R(\underline{e})(\Delta \mathbf{u} + \Delta \mathbf{d} + \Delta \mathbf{w}), \quad \Delta \mathbf{v}(0) = \Delta \mathbf{v}_0, \\ \Delta \dot{\mathbf{d}}(t) &= \Delta \mathbf{s}(t), \quad \Delta \mathbf{d}(0) = \Delta \mathbf{d}_0 \end{aligned} \quad (8)$$

Accelerometer dynamics is essentially due to anti-aliasing filter characterizing the f^2 shape noise (Canuto and Massotti, 2009) and transmission delay. Accelerometer error \mathbf{e}_{ak} is the combination of bias \mathbf{b}_{ak} , drift \mathbf{d}_{ak} , white noise (including quantization) \mathbf{w}_{ak} , and high-frequency f^2 -proportional noise \mathbf{h}_{ak} :

$$\mathbf{e}_{ak}(t) = \mathbf{b}_{ak}(t) + \mathbf{d}_{ak}(t) + \mathbf{w}_{ak}(t) + \mathbf{h}_{ak}(t) \quad (9)$$

A simplified model just accounts for delay and neglect high-frequency noise, because of anti-aliasing filter

$$\begin{aligned} \mathbf{y}_{ak}(t) &= \mathbf{a}_k(t - \tau_a) + \mathbf{e}_{ak}(t - \tau_a) \\ \mathbf{a}_k(t) &= \mathbf{u}_k(t) + \mathbf{d}_k(t) + \mathbf{w}_k(t), \end{aligned} \quad (10)$$

where \mathbf{a}_k is the total non gravitational acceleration in body coordinates. Relative position and rate may be obtained from GPS receivers and inter-satellite radio

transmission, and, in parallel, by the on-board optical metrology. Only the former is considered here. Let us denote GPS range and rate measurements as

$$\begin{aligned} \mathbf{y}_{rk}(t) &= \mathbf{r}_k(t) + \mathbf{e}_{rk}(t) \\ \mathbf{y}_{vk}(t) &= \mathbf{v}_k(t) + \mathbf{e}_{vk}(t) \end{aligned} \quad (11)$$

Formation measurements are obtained from (11) through LORF-to-inertial matrix $R(\mathbf{q}_o)$ and LORF rate $\boldsymbol{\omega}_o$

$$\begin{aligned} \Delta \mathbf{y}_r(t) &= R^T(\mathbf{q}_o)(\mathbf{y}_{r0} - \mathbf{y}_{r1})(t) \\ \Delta \mathbf{y}_v(t) &= R^T(\mathbf{q}_o)(\mathbf{y}_{v0} - \mathbf{y}_{v1})(t) - \boldsymbol{\omega}_o \times (\mathbf{y}_{r0} - \mathbf{y}_{r1})(t) \end{aligned} \quad (12)$$

Considering electrical propulsion for thrust range and lifetime issues, the thruster dynamics can be seen as a combination of flow dynamics (slow) and beam current dynamics (fast rise time < 0.1 s). The latter dominates close to drag-free and attitude bandwidth, below 1 Hz, and therefore thrust-to-force and torque relations may be accounted as static. Splitting the force/torque vector into three components, the static relations result in

$$\begin{bmatrix} m_k \mathbf{u}_{xk} \\ m_k \mathbf{u}_{hk} \\ J_k \mathbf{u}_{qk} \end{bmatrix} (t) = \begin{bmatrix} F_{xk} \\ \mathbf{F}_{hk} \\ \mathbf{M}_k \end{bmatrix} (t) = \begin{bmatrix} b_{xk} & \cong 0 \\ \mathbf{b}_{hk} & B_{hk} \\ \Delta \mathbf{b}_{qk} & B_{qk} \end{bmatrix} \begin{bmatrix} u_{mk} + d_{mk} \\ \mathbf{u}_{ik} + \mathbf{d}_{ik} \end{bmatrix} (t), \quad (13)$$

where \mathbf{u}_k in (6) has been split into along-track u_{xk} and cross-track & radial \mathbf{u}_{hk} , $\Delta \mathbf{b}_{qk}$ is due misalignment, and \mathbf{b}_{hk} is due to mini-thruster nominal inclination (as in Fig. 1). Note that u_{mk}, \mathbf{u}_{ik} denote mini and micro commanded thrusts affected by noise d_{mk}, \mathbf{d}_{ik} respectively, the latter being components of \mathbf{d}_k in Eq. (6) and of \mathbf{a}_{qk} in Eq. (3).

3 Requirements and control design

Control requirements split into

1. drag-free requirement from (10) and (6)

$$\mathbf{a}_k(t) = 0, \quad k = 0, 1, \quad (14)$$

2. attitude requirement from (3) and (4)

$$\begin{aligned} \mathbf{a}_{qk}(t) &= \underline{\mathbf{a}}(t) = \mathbf{a}_o(t) \\ \boldsymbol{\omega}_{qk}(t) &= \underline{\boldsymbol{\omega}}(t) = \boldsymbol{\omega}_o(t), \\ \mathbf{q}_k(t) &= \underline{\mathbf{q}}(t) = \mathbf{q}_o(t) \end{aligned} \quad (15)$$

3. formation requirement

$$\Delta \mathbf{r}(t) = \Delta \underline{\mathbf{r}} = [d \quad 0 \quad 0]^T. \quad (16)$$

Actual requirements admit residuals, which are expressed through spectral density bounds in case of attitude and drag-free variables, the latter being applicable to non gravitational CoM and angular accelerations. Formation residuals must be bounded by a box defined as

$$|\Delta r_j(t) - \Delta r_j| = \delta r_{j,\max}, \quad j = x, y, z. \quad (17)$$

Requirements must be completed with thrust bounds imposed by technology, as previously addressed in the introduction. Specifically

$$\begin{aligned} 0 < u_{\min} \leq u_m(t) \leq u_{\max} \\ 0 < u_{t,\min} \leq \max_t |\mathbf{u}_t(t)|_{\infty} \leq u_{t,\max} \end{aligned} \quad (18)$$

Any overshoot of the computed command is managed by control strategy, in order to not destroy the drag-free flight conditions, thus jeopardizing science. Propellant optimization may be added as a further objective (see Canuto, 2008).

Drag-free control may be designed as a pure disturbance rejection, with the constraints that accelerometer drift and bias are automatically rejected. The formation command \mathbf{u}_{fk} can be seen as follows

$$\mathbf{u}_k(t) = -\mathbf{d}_k(t) - \mathbf{b}_{ak}(t) - \mathbf{d}_{ak}(t) + \mathbf{u}_{fk}(t), \quad (19)$$

and the corresponding acceleration residuals hold (from (6) and (10))

$$\mathbf{a}_k(t) = \mathbf{w}_k(t) - \mathbf{b}_{ak}(t) - \mathbf{d}_{ak}(t) - \mathbf{w}_{ak}(t) - \mathbf{h}_{ak}(t). \quad (20)$$

Equations (19) and (20) clearly impose formation command and accelerometer noise to stay below drag-free bound. Drift in (20) is no detrimental as in a single spacecraft like GOCE, since the corresponding acceleration is bounded (drift is due basically to thermal fluctuations of the electronics) and largely lower than gravity, but it may destroy the formation in less than one day if not counteracted.

A detail design of formation control is presented in Canuto, Molano and Jimenez (2010). A generic formulation combines, in a multivariate law, drift cancellation and formation tracking so as to respect (17)

$$\Delta \mathbf{u}_f = \Delta \mathbf{b}_a + \Delta \mathbf{d}_a - K_r (\Delta \mathbf{r} - \Delta \mathbf{r}) - K_v \Delta \mathbf{v}, \quad (21)$$

where $\Delta \mathbf{u}_f = \mathbf{u}_{f0} - \mathbf{u}_{f1}$. A pair of challenging problems arise in implementing and tuning (21), as mentioned in Canuto, Molano and Jimenez (2010):

1. differential accelerometer drift and bias $\Delta \mathbf{b}_a + \Delta \mathbf{d}_a$ are not the only disturbance components in (8), since the main contribution comes from J_2 (Earth flattening) - the static component contributes to Hill dynamics together with $\boldsymbol{\omega}_O \times (\boldsymbol{\omega}_O \times \Delta \mathbf{r}_k + 2\Delta \mathbf{v}_k)$ -, and from the orbit eccentricity, entering $-\dot{\boldsymbol{\omega}}_O \times \Delta \mathbf{r}$. Such components must be fully excluded from (21) for both drag-free and thruster bounding reasons, which leads to a form of differential drag-free control.
2. thrust is minimized by an ad-hoc multivariate design of the feedback gains K_r, K_v , which exploits the cross-coupling properties of Hill's equation in a novel way (see Canuto, Molano and Jimenez, 2010).

Attitude control exploits both electric propulsion (micro-thrusters) and magnetic torquers. The resulting commanded acceleration has a similar form to (21) and holds, from (3) and (5),

$$\begin{aligned} \mathbf{u}_{qk}(t) &= \mathbf{a}_q + 0.5K_{qk} \mathbf{e}_k(t) / \underline{\rho}_k + K_{ok} (\boldsymbol{\omega} - \boldsymbol{\omega}_k) + \\ &+ (I - P_m) (J_k^{-1} (\boldsymbol{\omega}_k(t) \times J_k \boldsymbol{\omega}_k(t)) - \mathbf{a}_{qk}(t)) \quad , \\ \mathbf{u}_{mk}(t) &= P_m (J_k^{-1} (\boldsymbol{\omega}_k(t) \times J_k \boldsymbol{\omega}_k(t)) - \mathbf{a}_{qk}(t)) \end{aligned} \quad (22)$$

where feedback gains K_{qk} , K_{ok} allow LORF tracking and stability, whereas P_m is the projection of the rejected disturbance vector (Silani and Lovera, 2005), including gyro torque, on the subspace orthogonal to the instantaneous Earth magnetic field \mathbf{b}_E . The projected torque is contrasted by magnetic-torquer commanded acceleration \mathbf{u}_{mk} . All state variables entering the control laws (19), (21) and (22) are obtained from appropriate noise estimators in discrete time domain, which constitute the Embedded Model with the subset (2), (3), (6), (8), (10) and (13). All previous control law have been proved such to guarantee performance and stability. Space constraints prevent formal demonstration in the present paper.

4 Simulated results

All simulate results were obtained under the worst expected environment conditions, dictated by the highest solar activity, likely to be met during a 6-year mission. Spectral densities of drag-free residuals and target bound are shown in Fig. 3: at a first glance, drag-free bound is not respected at lower frequencies due to resonance peaks at orbit frequency, 0.2 mHz, and 2nd harmonics (J2). Actually spectral bound does not apply to periodic components, which must be bounded in terms of RMS. Thus, except for the orbit harmonics, bound is largely respected.

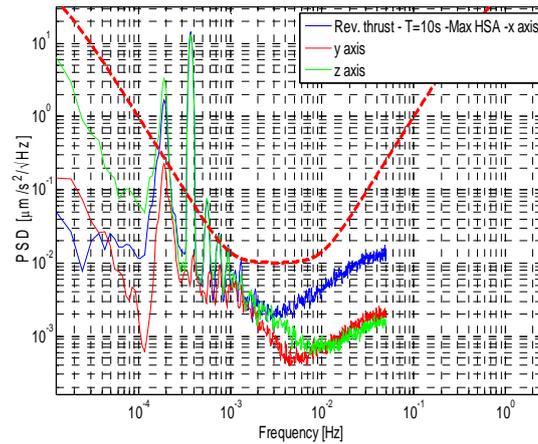


Fig. 3. Spectral density of drag-free residuals (single spacecraft).

3D formation tracking errors are shown in Fig. 4, which enlightens the long-term natural formation beat motion due to Earth flattening and eccentricity. It cannot be destroyed less a large increase in thrust peak. The beat motion carrier is the orbit period of about 5500 s. Note the peak of the radial motion being close to box limit in (17): 30 m versus 50 bound, entailing formation requirements were not loose.

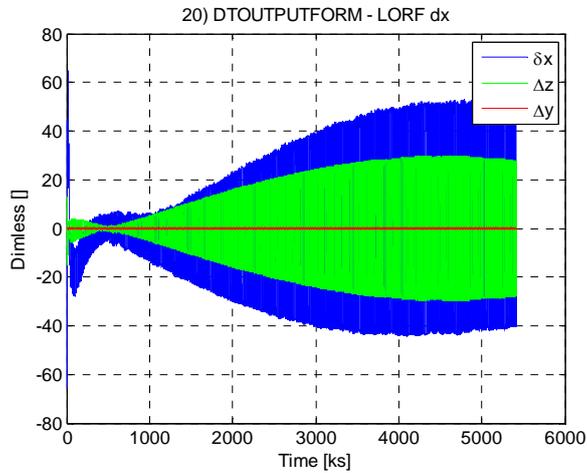


Fig. 4. Two-months formation relative position.

Table 1 reports thruster performance and target under the expected worst-case conditions. Note propellant mass and peak power are outside target, which underlines mission requirements being at a technology border.

No.	Type	Unit	Value	Target
0	Propellant mass	kg	70	50
1	Average power	W	460	500
2	Peak power	W	1250	1000

5 Conclusions and acknowledgments

An overview of the control challenges in view of a long-distance, drag-free, low-Earth-orbit spacecraft formation, together with an outline of their solution have been presented, supported by simulation results. Part of the work has been done under a grant of the European Space Agency to Politecnico di Torino in collaboration with Thales Alenia Space Italia, Turin, Italy.

References

1. Canuto, E. (2007) Embedded Model Control: outline of the theory. *ISA Trans.* **46** (3), 363-377.
2. Canuto, E (2008) Drag-free and attitude control for the GOCE satellite. *Automatica*, **44**, 1766-1780.
3. Canuto, E. and Massotti, L. (2009) All-propulsion design of the drag-free and attitude control of the European satellite GOCE. *Acta Astronautica*, **64**, 325-344.

4. Canuto, E, Molano, A. and Massotti, L. (2010) Drag-free control of the GOCE satellite: noise and observer design. *IEEE Trans. on Control System Technology*, **18**, 501-509.
5. Canuto, E., Molano, A., Jimenez, J. and Perez C. (2010) Long-distance, drag-free, low-thrust, LEO formation control. Submitted to *Chinese Control Conference 2010 (CCC)*, Beijing, July 29-31, 2010.
6. Inalhan, G., Tillerson, M. and How, J.P. (2002) Relative dynamics and control of spacecraft formations in eccentric orbits. *J. Guidance, Control and Dynamics*, **25** (1), 48-60.
7. Loeb, H.W., Scharner, K.-H., Weiss, St., Feili, D. and Meyer, B.K. (2004) Development of RIT-microthrusters. *Proc. 55th Int. Astronautical Congr.*, Vancouver, Canada.
8. Silani, E., and Lovera, M. (2005) Magnetic spacecraft attitude control: A survey and some new results, *Control Engineering Practice, Special Section on Aerospace Control*, **13** (3), 357-371.

The Flexible Dirichlet family: some inferential issues

S. Migliorati¹, G. S. Monti¹, and A. Ongaro¹

¹ Department of Statistics, University of Milano-Bicocca,
Via Bicocca degli Arcimboldi 8, 20126 Milano, Italy.
E-mail for correspondence: sonia.migliorati@unimib.it

Abstract. The Flexible Dirichlet distribution (Ongaro et al., 2008, [5]) has been recently introduced to model compositional data. It is a generalization of the Dirichlet which preserves some of its good mathematical properties and, at the same time, exhibits a richer dependence structure which allows various forms of dependence relevant for compositional data, independence cases being identified by suitable parameter configurations.

Here we investigate the nature of the dependence introduced by the new distribution. Furthermore we develop suitable likelihood-based testing procedures to assess the presence of dependence relations of particular impact in applications. Their performances will be evaluated by means of Monte Carlo experiments.

Keywords: Generalizations of Dirichlet distribution, Finite mixture, Compositional data, Neutrality, Likelihood.

1 Introduction

In many problems data consist of vectors of proportions, such as chemical constituents of a substance, and are, therefore, subject to a unit sum constraint. This type of data, called compositional, arise naturally in a great variety of disciplines such as archeology, biology, economics, environmetrics, psephology, medicine, psychology, etc..

The most well known distribution for compositional data is the Dirichlet which possesses several good statistical and mathematical properties, such as closure under amalgamation and subcomposition, as well as easiness of parameter interpretation. However it is only suitable for modeling data exhibiting the maximum degree of independence compatible with compositions.

The Flexible Dirichlet (FD) distribution (Ongaro et al., 2008, [5]) allows to overcome such serious drawback by accounting for various types of dependence.

After reviewing some properties of such distribution (Sections 2 and 3) we focus on its (in)dependence structure. On the one hand we study the type of non neutrality provided for by the model by analyzing the influence of a given subset of variables on the subcomposition formed by the other ones (Section 4). On the other we develop suitable testing procedures to assess the presence of independence relations of particular impact in applications (Section 5).

2 The Flexible Dirichlet distribution

The Dirichlet distribution $\underline{X} \sim \mathcal{D}^D(\underline{\alpha})$, with $\underline{\alpha} = (\alpha_1, \dots, \alpha_D) \in \mathbb{R}_+^D$, takes values on the unitary simplex $S^D = \{\underline{x} : x_i \geq 0, i = 1, \dots, D \text{ and } \sum_{i=1}^D x_i = 1\}$.

Such distribution can be obtained by normalizing a vector (basis) of independent, equally scaled Gamma random variables (r.v.s) and it is closed under operations of marginalization, conditioning, amalgamation and sub-composition, the consequent distributions being simply related to the full one.

The FD distribution is achieved by normalizing a basis of dependent r.v.s which contains equally scaled Gamma independent variates as a particular case. Let $W_i \sim Ga(\alpha_i)$ ($\alpha_i > 0$) denote such Gamma r.v.s ($i = 1, \dots, D$) and let $U \sim Ga(\tau)$ ($\tau > 0$) denote a further independent Gamma r.v. which is allocated to the i^{th} component of the basis with probability p_i ($0 < p_i < 1$ and $\sum_{i=1}^D p_i = 1$). Then, the new basis $\underline{Y} = (Y_1, \dots, Y_D)$ is defined as $Y_i = W_i + Z_i U$, $i = 1, \dots, D$, where $\underline{Z} = (Z_1, \dots, Z_D)$ is a multinomial vector independent from U and from the W_i 's which is equal to \underline{e}_i with probability p_i where \underline{e}_i is a vector of zeros except for the i^{th} element which is one.

The normalized vector $\underline{X} = (\frac{Y_1}{Y^+}, \dots, \frac{Y_D}{Y^+})$, (where $Y^+ = \sum_{i=1}^D Y_i$), has a FD distribution denoted by $FD^D(\underline{\alpha}, \underline{p}, \tau)$ and it is a finite mixture of Dirichlet distributions:

$$FD^D(\underline{\alpha}, \underline{p}, \tau) = \sum_{i=1}^D p_i \mathcal{D}^D(\underline{\alpha} + \tau \underline{e}_i). \quad (1)$$

Therefore, its density function can be expressed as

$$f_{FD}(\underline{x}; \underline{\alpha}, \underline{p}, \tau) = \frac{\Gamma(\alpha^+ + \tau)}{\prod_{r=1}^D \Gamma(\alpha_r)} \left(\prod_{r=1}^D x_r^{\alpha_r - 1} \right) \sum_{i=1}^D p_i \frac{\Gamma(\alpha_i)}{\Gamma(\alpha_i + \tau)} x_i^\tau \quad (2)$$

where \underline{x} belongs to the unitary simplex and $\alpha^+ = \sum_{i=1}^D \alpha_i$.

Here we have adopted a slightly smaller parameter space (the interior of the original one) than in Ongaro et al. 2008 [5] which is more tractable from a mathematical point of view without losing in generality in terms of independence relations.

For completeness we now report some useful properties of the FD suitably adjusted to the new parameter space.

The FD includes the Dirichlet as an inner point: $FD^D(\underline{\alpha}, \underline{p}, \tau) \equiv \mathcal{D}^D(\underline{\alpha})$ if and only if $\tau = 1$ and $p_i = \alpha_i / \alpha^+$, $\forall i = 1, \dots, D$.

The first two moments can be expressed as:

$$E(X_i) = \frac{\alpha_i + p_i \tau}{\alpha^+ + \tau}$$

$$Var(X_i) = \frac{E(X_i)(1 - E(X_i))}{(\alpha^+ + \tau + 1)} + \frac{\tau^2 p_i (1 - p_i)}{(\alpha^+ + \tau)(\alpha^+ + \tau + 1)}$$

$$Cov(X_i, X_r) = -\frac{E(X_i)E(X_r)}{(\alpha^+ + \tau + 1)} - \frac{\tau^2 p_i p_r}{(\alpha^+ + \tau)(\alpha^+ + \tau + 1)}.$$

Thus, unlike the Dirichlet, the FD distribution accounts for components with the same mean but different variances or for covariances which do not show proportionality with respect to the product of means.

In order to characterize marginal and conditional distributions, it is useful to adopt the following notation. Given a partition (of order 1) $\underline{X} = (X_1, \dots, X_k | X_{k+1}, \dots, X_D) = (\underline{X}_1, \underline{X}_2)$ we shall denote the corresponding totals by $X_1^+ = \sum_{i=1}^k X_i$ and $X_2^+ = \sum_{i=k+1}^D X_i$ and, in an analogous way, we shall define the quantities $\underline{\alpha}_1, \underline{\alpha}_2, \alpha_1^+, \alpha_2^+, p_1, p_2, p_1^+$ and p_2^+ . Moreover we shall indicate the two subcompositions by $\underline{S}_1 = \frac{(X_1, \dots, X_k)}{X_1^+}$ and $\underline{S}_2 = \frac{(X_{k+1}, \dots, X_D)}{X_2^+}$ and the amalgamation (vector of totals) by $\underline{T} = (X_1^+, X_2^+)$.

First of all, the FD distribution is closed under marginalization, i.e.:

$$(\underline{X}_1, 1 - X_1^+) \sim FD^{k+1}(\underline{\alpha}_1, \alpha^+ - \alpha_1^+, p_1, 1 - p_1^+, \tau). \tag{3}$$

Furthermore, its (normalized) conditional distributions are mixtures of a FD and of a Dirichlet. More precisely:

$$\frac{\underline{X}_1}{1 - x_2^+} | \underline{X}_2 = \underline{x}_2 \sim \underline{S}_1 | \underline{X}_2 = \underline{x}_2$$

has distribution:

$$p(\underline{x}_2)FD^k\left(\underline{\alpha}_1, \frac{p_1}{p_1^+}, \tau\right) + (1 - p(\underline{x}_2))\mathcal{D}^k(\underline{\alpha}_1) \tag{4}$$

where

$$p(\underline{x}_2) = \frac{p_1^+}{p_1^+ + q(\underline{x}_2)} \tag{5}$$

and

$$q(\underline{x}_2) = \frac{\Gamma(\alpha_1^+ + \tau)}{\Gamma(\alpha_1^+)(1 - x_2^+)^\tau} \sum_{i=k+1}^D p_i \frac{\Gamma(\alpha_i)}{\Gamma(\alpha_i + \tau)} x_i^\tau. \tag{6}$$

The FD is also closed under permutation (the parameters of the permuted random vector being simply the permutation of the original parameters) and under amalgamation (where the parameters of the amalgamation can be easily obtained by summing up the α_i 's and p_i 's within each group of the partition).

Finally, the distribution of subcompositions from a FD can be easily derived. For example: $\underline{S}_1 \sim p_1^+ FD^k\left(\underline{\alpha}_1, \frac{p_1}{p_1^+}, \tau\right) + (1 - p_1^+) \mathcal{D}^k(\underline{\alpha}_1)$.

Notice also that properties concerning amalgamation and subcompositions do hold for partitions of any order (i.e. into an arbitrary number of subsets).

3 Independence relationships

Clearly the components of a random vector defined on the simplex cannot be independent because of the unit-sum constraint. That is why a large variety of ad hoc forms of independence has been developed in the literature (see for example Aitchison, 1980, [1], and 2003, [2]), most of which can be expressed in terms of subcompositions and amalgamation. Focusing on partitions of order 1 for the sake of simplicity, we shall mainly consider partition independence ($\underline{S}_1 \perp \underline{S}_2 \perp \underline{T}$, where \perp stands for independence), neutrality on the right ($\underline{S}_2 \perp (\underline{S}_1, \underline{T})$) and neutrality on the left ($\underline{S}_1 \perp (\underline{S}_2, \underline{T})$).

The Dirichlet distribution can be shown to possess all the above independence properties and it can be properly considered as the model of maximum independence compatible with unit-sum constrained r.v.s.

Vice versa the FD exhibits a rich dependence structure, various forms of independence corresponding to suitable parameter configurations. Let us focus on the following ones which will prove to be the most interesting:

1. $\tau = 1$ and $p_i = \alpha_i/\alpha^+$, ($i = 1, \dots, D$), i.e. $\underline{X} \sim \mathcal{D}^D(\alpha)$;
2. $\tau = 1$ and $\frac{p_i}{p_1^+} = \frac{\alpha_i}{\alpha_1^+}$, ($i = 1, \dots, k$);
3. $\tau = 1$ and $\frac{p_i}{p_2^+} = \frac{\alpha_i}{\alpha_2^+}$, ($i = k + 1, \dots, D$).

It can be proved that the $FD^D(\underline{\alpha}, \underline{p}, \tau)$ is neutral on the left, i.e. $\underline{S}_1 \perp (\underline{S}_2, \underline{T})$, if and only if either condition **1.** or condition **2.** is satisfied. Analogously, conditions for right neutrality can be obtained: we have $\underline{S}_2 \perp (\underline{S}_1, \underline{T})$ if and only if either condition **1.** or condition **3.** is satisfied. Furthermore, the $FD^D(\underline{\alpha}, \underline{p}, \tau)$ shows partition independence, i.e. $\underline{S}_1 \perp \underline{S}_2 \perp \underline{T}$, if and only if either condition **1.** or both condition **2.** and **3.** are satisfied.

Finally, it is noticeable that conditions for independence relations to hold can be generalized to higher order partitions. For example, a partition of order 2 shows partition independence if and only if either condition **1.** is satisfied or $\tau = 1$ and in at least two of the three subsets the α_i 's are proportional to the p_i 's.

4 Dependence pattern

Whenever in a given model a type of independence is absent, it is of statistical interest to analyze the form of the consequent dependence. This gives rise to a number of relationships of potential importance. Here we shall focus on the independence concept of neutrality whose relevance and generality, both from a theoretical and an applicative perspective, clearly emerges from the literature. Such concept, first introduced by Connor and Mosimann, 1969 [3], has to do with the consequences of eliminating a certain number of components on the relative proportions of the remaining ones (i.e. the corresponding

subcomposition). To illustrate the concept, suppose that the researcher’s interest is on (X_1, \dots, X_{D-1}) ; then X_D is said neutral, and it can therefore be neglected, only if it has no influence on $(X_1/(1 - X_D), \dots, X_{D-1}/(1 - X_D))$. For example, consider a household budget analysis where the total expenditures are classified into a number of commodity categories. Then it may be of interest to understand whether, for instance, the amount spent on foodstuffs affects the expenditure pattern (subcomposition) of the other categories.

More generally, using our notation, a vector \underline{X}_2 is neutral if it is independent of \underline{S}_1 . It is important to observe that such notion coincides with the above introduced neutrality on the left due to the one-to-one correspondence between \underline{X}_2 and $(\underline{S}_2, \underline{T})$.

If \underline{X}_2 is not neutral, then it is of obvious interest to analyze how it does affect the composition of the remaining variables. Such issue can be explored within the FD model by considering the conditional mean effect.

Proposition 1

Let $\underline{S}_1 = (S_{11}, \dots, S_{1k})$, then for $i = 1, \dots, k$

$$E(S_{1i} | \underline{X}_2 = x_2) = (1 - p(x_2)w) \frac{\alpha_i}{\alpha_1^+} + p(x_2)w \frac{p_i}{p_1^+} \tag{7}$$

where $p(x_2)$ is given by (5) and $w = \tau/(\tau + \alpha_1^+)$.

Proof

The result follows, after some algebraic work, from (4) and knowledge of the first moment of the FD. □

The conditional mean is easily seen monotone in $p(x_2)$ and therefore in each x_j , ($j = k+1, \dots, D$), being $p(x_2)$ decreasing in each x_j . More precisely, it varies from $(1 - w) \frac{\alpha_i}{\alpha_1^+} + w \frac{p_i}{p_1^+}$ when $x_2^+ = 0$ to $\frac{\alpha_i}{\alpha_1^+}$ when $x_2^+ \rightarrow 1$. In particular, it is increasing (decreasing) when $\frac{p_i}{p_1^+} < \frac{\alpha_i}{\alpha_1^+}$ ($\frac{p_i}{p_1^+} > \frac{\alpha_i}{\alpha_1^+}$) and it is constant when $\frac{p_i}{p_1^+} = \frac{\alpha_i}{\alpha_1^+}$, thus making simple to model both positive and negative dependences.

The parameter τ determines the range of variation of the conditional mean, the bigger τ the larger such range. Particularly simple expressions are obtained when $\tau = 1$; if moreover p_i is proportional to α_i for $i = k+1, \dots, D$, then the conditional mean depends on \underline{x}_2 only through the sum x_2^+ .

5 Testing independence

A convenient strategy to analyze the various forms of independence is to order them from the strongest to the weakest and then to test them in such order, proceeding to the next level only in case of rejection of the preceding one.

In general the first hypothesis to be tested is the Dirichlet model one as it implies any other independence. Then, given a partition of order 1 of interest, one can test partition independence first and finally, at the same level, neutrality on the left and/or on the right. Notice that the last two properties are equivalent to partition independence.

Any of the above hypotheses can be tested through a suitable likelihood ratio test with asymptotic chi-square distribution.

Obviously there is no guarantee of a complete coherence among decisions taken at the various steps: for example rejection of partition independence may occur without rejecting neither neutrality on the left nor neutrality on the right. If such coherence is thought essential then one might look for alternative strategies such as intersection-union tests: reject partition independence iff at least one of the two neutrality likelihood ratio tests rejects. In this case, to obtain a level α test for partition independence a level $\alpha/2$ can be adopted for the other two tests. Anyway, notice that such solution leads to a conservative test.

Clearly, if the researcher is interested in just one particular type of independence, i.e. neutrality, she/he does not need to follow the above scheme.

The construction of the above mentioned likelihood ratio tests requires the unconstrained maximization of the likelihood as well as the constrained maximization under the various hypotheses, which are critical issues given the mixture structure of the model. The former problem has been tackled in Migliorati et al. 2008 [4] where an E-M- algorithm has been adopted with initial values obtained by combining the k -mean clustering algorithm for estimating the p_i 's and a two step method of moments for τ and $\underline{\alpha}$.

Let us now focus on the issues arising from constrained maximization. First let us consider the null hypothesis relative to condition **1.**, i.e. $H_0 : \underline{X} \sim \mathcal{D}^D(\alpha)$. The maximization of the likelihood under H_0 has been performed by applying the Newton-Raphson algorithm (Ronning,1989 [6]) using the method of moments to obtain the starting values. The test statistic distribution has been approximated by a chi-square with D degrees of freedom according to Wilks' theorem.

The other hypotheses require a more complex procedure. Maximization of the likelihood under such null hypotheses is best achieved by constructing suitable profile likelihoods which exploit specific factorization properties of the FD model. In particular, under the null hypothesis of partition independence

$$H_0 : \tau = 1; \frac{p_i}{p_1^+} = \frac{\alpha_i}{\alpha_1^+}, (i = 1, \dots, k); \frac{p_i}{p_2^+} = \frac{\alpha_i}{\alpha_2^+}, (i = k + 1, \dots, D)$$

the distribution of \underline{X} can be conveniently represented through the distribution of \underline{S}_1 , \underline{S}_2 and \underline{T} , which are independent with $\underline{S}_1 \sim \mathcal{D}^k(\underline{\alpha}_1)$, $\underline{S}_2 \sim \mathcal{D}^{D-k}(\underline{\alpha}_2)$ and $\underline{T} \sim FD^2(\alpha_1^+, \alpha_2^+, p_1^+, p_2^+, \tau = 1)$ (for a proof see Section 5 of Ongaro et al. 2008 [5]). Such formulation has the advantage of automatically incorporating the null hypothesis constraints. Furthermore it suggests

to consider the profile likelihood for α_1^+ and α_2^+ as only such parameters appear in more than one of the above distributions. The profile can be easily constructed by separate maximization of the likelihood relative to the three distributions. The test statistic distribution has been approximated by a chi-square with $D - 1$ degrees of freedom.

The null hypothesis of left neutrality:

$$H_0 : \tau = 1; \frac{p_i}{p_1^+} = \frac{\alpha_i}{\alpha_1^+}, (i = 1, \dots, k)$$

can be dealt with a similar method: \underline{X} is best represented through the distribution of \underline{S}_1 and \underline{X}_2 which are independent with $S_1 \sim \mathcal{D}^k(\underline{\alpha}_1)$ and $\underline{X}_2 \sim FD^{D-k}(\alpha_{k+1}, \dots, \alpha_D, \alpha_1^+, p_{k+1}, \dots, p_D, p_1^+, \tau = 1)$. This leads to consider the profile for α_1^+ which is the only parameter shared by the two distributions. The test statistic distribution has been approximated by a chi-square with k degrees of freedom. Clearly the hypothesis of right neutrality can be tested in an analogous way.

A simulation study of the performances of the above tests has been carried out with 10,000 replications for different values of n and of the parameter vector. The following tables report the simulated (real) significance levels against a 5% nominal one and some values of the simulated power.

Table 1 refers to $H_0 : \underline{X} \sim \mathcal{D}^D(\underline{\alpha})$ and takes into consideration the following parameter configurations where the null hypothesis is true only in cases (a), (b) and (c):

- (a) $\underline{X} \sim D^3(\underline{\alpha} = (1, 1, 1))$
- (b) $\underline{X} \sim D^4(\underline{\alpha} = (0.5, 0.3, 0.7, 0.6))$
- (c) $\underline{X} \sim D^5(\underline{\alpha} = (6, 4, 3, 1, 8))$
- (d) $\underline{X} \sim FD^3(\underline{\alpha} = (1, 1, 1), \underline{p} = (0.45, 0.25, 0.3), \tau = 2)$
- (e) $\underline{X} \sim FD^4(\underline{\alpha} = (0.6, 0.3, 0.5, 0.2), \underline{p} = (0.15, 0.35, 0.3, 0.2), \tau = 4)$
- (f) $\underline{X} \sim FD^5(\underline{\alpha} = (8, 3, 4, 2, 10), \underline{p} = (0.15, .35, 0.15, 0.2, 0.15), \tau = 6)$

Table 1. Proportion of rejections at 5% level.

case	$n = 50$	$n = 100$	$n = 300$
(a)	0.058	0.061	0.052
(b)	0.047	0.05	0.048
(c)	0.054	0.046	0.047
(d)	0.217	0.336	0.757
(e)	0.996	0.998	1
(f)	0.991	1	1

The first three rows highlight a good performance of the simulated significance level for all models considered and all sample sizes. It is also noticeable

that the power quickly converges to 1 for increasing sample sizes except for the case (d) whose parameter configuration is however quite close to the null.

For lack of space we report simulations only for the neutrality case as it is computationally more demanding than the partition independence one. Table 2 reports the simulation results referred to the left neutrality hypothesis where the original composition has been partitioned as $(X_1, X_2, X_3|X_4, X_5)$. We considered the following parameter configurations:

- (a) $\underline{X} \sim FD^5(\underline{\alpha} = (6, 5, 13, 10, 6), \underline{p} = (0.4 \cdot (6, 5, 13)/24, 0.3, 0.3), \tau = 1)$
- (b) $\underline{X} \sim FD^5(\underline{\alpha} = (5, 10, 20, 6, 9), \underline{p} = (0.5 \cdot (5, 10, 20)/35, 0.4, 0.1), \tau = 1)$
- (c) $\underline{X} \sim FD^5(\underline{\alpha} = (0.5, 2, 5, 0.6, 1), \underline{p} = (0.25, 0.3, 0.2, 0.1, 0.15), \tau = 2)$
- (d) $\underline{X} \sim FD^5(\underline{\alpha} = (5, 10, 20, 6, 9), \underline{p} = (0.3, 0.2, 0.1, 0.1, 0.3), \tau = 3)$

where the null hypothesis is true only in cases (a) and (b).

Table 2. Proportion of rejections at 5% level.

case	$n = 300$	$n = 500$	$n = 1000$
(a)	0.101	0.075	0.064
(b)	0.109	0.069	0.062
(c)	0.468	0.695	0.966
(d)	0.198	0.215	0.366

The performance of the simulated significance level appears to be satisfactory even though the convergence is slower than in the case of Table 1. Concerning the power of the test, some convergence difficulties emerge in case (d) which deserves further investigation.

References

1. Aitchison, J., “The Statistical Analysis of Compositional Data (with discussion)”, *Journal of the Royal Statistical Society* 44, 139–177 (1982).
2. Aitchison, J. *The Statistical Analysis of Compositional Data*, The Blackburn Press, London (2003).
3. Connor, J.R., and Mosimann, J.E., “Concepts of Independence for Proportions with a Generalization of the Dirichlet Distribution”, *Journal of the American Statistical Association* 64, 194–206 (1969).
4. Migliorati, S., Monti, G.S., and Ongaro, A., “E–M Algorithm: an Application to a Mixture Model for Compositional Data”, *Proceedings of the 44th Scientific Meeting of the Italian Statistical Society* (2008).
5. Ongaro, A., Migliorati, S., Monti, G.S., “A new distribution on the simplex containing the Dirichlet family”, *Proceedings of CODAWORK’08, The 3rd Compositional Data Analysis Workshop, University of Girona, Spain* (2008).
6. Ronning, “Maximum likelihood estimation of Dirichlet distributions”, *Journal of Statistical Computation and Simulation* 32, 215–221 (1989).

Compound Binomial Risk Model [★]

Leda D.Minkova

Faculty of Mathematics and Informatics
Sofia University "St. Kl.Ohridski",
Sofia, Bulgaria
(e-mail: leda@fmi.uni-sofia.bg)

Abstract. In this note the counting process in the insurance risk model is a compound Binomial process. The particular case of geometric compounding distribution is analyzed. The counting process is called Inflated-parameter binomial process (I - Binomial process). Some basic properties are given. The corresponding risk model is called I - Binomial risk model. The joint probability distribution of the time to ruin and the deficit after ruin occurs is studied. The case of exponentially distributed claims is given.

Keywords: Discrete distributions,, inflated - parameter distributions,, compound distribution.

1 Introduction

Consider the standard risk model $\{X(t), t \geq 0\}$, defined on the complete probability space (Ω, \mathcal{F}, P) and given by

$$X(t) = ct - \sum_{k=1}^{N(t)} Z_k, \quad \left(\sum_1^0 = 0 \right). \quad (1)$$

Here c is a positive real constant representing the risk premium rate. The sequence $\{Z_k\}_{k=1}^{\infty}$ of non-negative mutually independent identically distributed random variables is independent of the counting process $N(t)$, $t \geq 0$. The claim sizes $\{Z_k\}_{k=1}^{\infty}$ are distributed as the random variable Z with distribution function F , $F(0) = 0$ and mean value $\mu = EZ < \infty$.

In this paper we suppose that the counting process $N(t)$ has a compound binomial distribution, see [1]. Suppose that $N(t) = \sum_{i=1}^{N_1(t)} Y_i$, where Y_1, Y_2, \dots are independent identically distributed random variables, independent of $N_1(t)$ and for $\alpha > 0$ and $n \geq 1$, $N_1(t) \sim Bi(n, \frac{t}{\alpha})$. Let Y denote the compounding random variable. Here we suppose that for $\rho \in [0, 1)$, $Y \sim Ge_1(1 - \rho)$ with probability mass function

$$P(Y = m) = (1 - \rho)\rho^{m-1}, \quad m = 1, 2, \dots$$

The compound binomial process $N(t)$ has Inflated - parameter binomial distribution ([3] and [5]) and is called I-Binomial process.

[★] This paper is partially supported by Sofia University grant 028/2009

In this paper the counting process $N(t)$ is defined as a birth process. Some properties are given and the application in insurance risk model is analyzed. We consider the particular case of exponentially distributed claims.

2 I-Binomial process

The I-Binomial process as a generalized birth process is defined in [4]. The transition probabilities are given by the following postulates:

$$P(N(t+h) = n \mid N(t) = m) = \begin{cases} 1 - (1-\rho) \frac{n\alpha}{(\alpha-t)^2} \sum_{k=1}^{\infty} \left[1 - (1-\rho) \frac{\alpha}{\alpha-t}\right]^{k-1} h + o(h), & n = m, \\ (1-\rho) \frac{n\alpha}{(\alpha-t)^2} \left[1 - (1-\rho) \frac{\alpha}{\alpha-t}\right]^{k-1} h + o(h), & n = m+k, k = 1, 2, \dots, \end{cases}$$

for every $m = 0, 1, \dots$, where $o(h) \rightarrow 0$ as $h \rightarrow 0$.

If $P_m(t) = P(N(t) = m)$, $m = 0, 1, 2, \dots$, the above postulates yield the following Kolmogorov forward equations:

$$P'_0(t) = -\frac{n}{\alpha-t} P_0(t),$$

$$P'_m(t) = -\frac{n}{\alpha-t} P_m(t) + (1-\rho) \frac{\alpha}{n} \left(\frac{\alpha}{\alpha-t}\right)^2 \sum_{k=1}^m \left[1 - (1-\rho) \frac{\alpha}{\alpha-t}\right]^{k-1} P_{m-k}(t), \tag{2}$$

for $m = 1, 2, \dots$. The solution of (2) with conditions

$$P_0(0) = 1 \quad \text{and} \quad P_m(0) = 0, \quad m = 1, 2, \dots$$

is given by

$$P(N_t = m) = \begin{cases} \left(1 - \frac{t}{\alpha}\right)^n, & m = 0 \\ \sum_{i=1}^{m \wedge n} \binom{n}{i} \binom{m-1}{i-1} \left[(1-\rho) \frac{t}{\alpha}\right]^i \left(1 - \frac{t}{\alpha}\right)^{n-i} \rho^{m-i}, & m = 1, 2, \dots \end{cases} \tag{3}$$

This is just the Inflated - parameter binomial distribution with parameters $\frac{t}{\alpha}$, ρ and n , say $\text{IBi}(\frac{t}{\alpha}, \rho, n)$ (see [3] and [5]). In the case of $\rho = 0$ (3) coincides with the usual binomial distribution.

This leads to the following definition

Definition 1 *The counting process $\{N(t), t \geq 0\}$ is said to be I - Binomial process, if it starts at zero, $N(0) = 0$ and for each $t > 0$, the distribution of $N(t)$ is given by (3).*

2.1 Properties of the I-Binomial process

Denote $S_m = T_1 + T_2 + \dots + T_m$, $m = 1, 2, \dots$, the waiting time until the m th event. One of the basic properties of the I-Binomial process is given in the next theorem.

Theorem 1 *Let $N(t)$ has the $IBi(\frac{t}{\alpha}, \rho, n)$ distribution (3). Then the waiting time until the m th event has the following probability density function (p.d.f.)*

$$f_{S_m}(t) = \frac{n}{\alpha} \sum_{i=0}^{(m-1) \wedge (n-1)} \binom{n-1}{i} \binom{m-1}{i} \left[(1-\rho) \frac{t}{\alpha} \right]^i \left(1 - \frac{t}{\alpha} \right)^{n-i-1} \rho^{m-i-1}. \tag{4}$$

2.2 Moments

The mean value and the variance of I - Binomial process are given by

$$EN(t) = \frac{nt}{(1-\rho)\alpha}$$

and

$$Var(N(t)) = \frac{n}{(1-\rho)^2} \left[1 - \frac{t}{\alpha} + \rho \right] \frac{t}{\alpha} = EN(t) \left[\frac{1+\rho}{1-\rho} - \frac{t}{(1-\rho)\alpha} \right]$$

3 Application to Risk Theory

We consider the risk model (1), where $N(t)$ is I - Binomial process and will call this process I - Binomial risk model.

The relative safety loading θ is defined by

$$\theta = \frac{c\alpha(1-\rho)}{n\mu} - 1,$$

and in the case of positive safety loading $\theta > 0$, $c > \frac{n\mu}{\alpha(1-\rho)}$.

We are interested in the probability that ruin occurs and the deficit at the time of ruin does not exceeds a given amount $y > 0$.

Let $\tau = \inf\{t : X(t) < -u\}$ with the convention of $\inf \emptyset = \infty$ be the time to ruin of an insurance company having initial capital $u \geq 0$. We denote by $\Psi(u) = P(\tau < \infty)$ the ruin probability and $\Phi(u) = 1 - \Psi(u)$ the nonruin probability.

In the following we use the notation of [2]. Let $G(u, y)$ be the joint probability distribution of the time to ruin τ and the deficit in prior to ruin $D = |U(\tau)|$ i.e.

$$G(u, y) = P(\tau < t, D \leq y). \tag{5}$$

and

$$\lim_{y \rightarrow \infty} G(u, y) = \Psi(u).$$

Using the postulates we have

$$\begin{aligned} G(u, y) &= \\ &= \left(1 - \frac{n}{\alpha-t}h\right) G(u+ch, y) + (1-\rho)\frac{n}{\alpha} \left(\frac{\alpha}{\alpha-t}\right)^2 h \sum_{k=1}^{\infty} \left[1 - (1-\rho)\frac{\alpha}{\alpha-t}\right]^{k-1} \\ &\times \left[\int_0^{u+ch} G(u+ch-x, y) dF^{*k}(x) + (F^{*k}(u+ch+y) - F^{*k}(u+ch)) \right] + o(h), \end{aligned}$$

where $F^{*k}(x)$, $k = 1, 2, \dots$ is the distribution function of $Z_1 + Z_2 + \dots + Z_k$. Rearranging the terms leads to

$$\begin{aligned} \frac{G(u+ch, y) - G(u, y)}{ch} &= \\ &= \frac{n}{(\alpha-t)c} G(u+ch, y) - (1-\rho)\frac{n}{\alpha c} \left(\frac{\alpha}{\alpha-t}\right)^2 \sum_{k=1}^{\infty} \left[1 - (1-\rho)\frac{\alpha}{\alpha-t}\right]^{k-1} \\ &\times \left[\int_0^{u+ch} G(u+ch-x, y) dF^{*k}(x) + (F^{*k}(u+ch+y) - F^{*k}(u+ch)) \right] + \frac{o(h)}{h}. \end{aligned}$$

Let

$$H(x) = (1-\rho)\frac{\alpha}{\alpha-t} \sum_{k=1}^{\infty} \left(1 - (1-\rho)\frac{\alpha}{\alpha-t}\right)^{k-1} F^{*k}(x) \quad (6)$$

be the non defective probability distribution function of the claims with

$$H(0) = 0, \quad H(\infty) = 1.$$

By letting $h \rightarrow 0$ we obtain the following differential equation

$$\frac{\partial G(u, y)}{\partial u} = \frac{n}{c(\alpha-t)} \left[G(u, y) - \int_0^u G(u-x, y) dH(x) - [H(u+y) - H(u)] \right]. \quad (7)$$

Theorem 2 *The function $G(0, y)$ is given by*

$$G(0, y) = \frac{n}{c(\alpha-t)} \int_0^y [1 - H(u)] du. \quad (8)$$

Proof. Integrating (7) from 0 to ∞ with $G(\infty, y) = 0$ leads to

$$\begin{aligned}
 & -G(0, y) = \\
 & = \frac{n}{c(\alpha - t)} \left[\int_0^\infty G(u, y) du - \int_0^\infty \int_0^u G(u - x, y) dH(x) du - \int_0^\infty (H(u + y) - H(u)) du \right]
 \end{aligned}$$

The change of variables in the double integral and simple calculations yield

$$G(0, y) = \frac{n}{c(\alpha - t)} \int_0^\infty [H(u + y) - H(u)] du$$

and (8). △

Theorem 3 *The ruin probability with $u = 0$ is given by*

$$\Psi(0) = \frac{n\mu}{(1 - \rho)\alpha c}. \tag{9}$$

Proof. According (8)

$$\Psi(0) = \lim_{y \rightarrow \infty} G(0, y) = \frac{n}{c(\alpha - t)} \int_0^\infty [1 - H(x)] dx.$$

Let X be a random variable with distribution function $H(x)$. By the definition of $H(x)$ and $EZ = \mu$ we obtain

$$EX = \frac{\mu(\alpha - t)}{(1 - \rho)\alpha}.$$

Using the fact that $EX = \int_0^\infty [1 - H(x)] dx$ we obtain (9). △

3.1 Exponentially distributed claims

Let us consider the case of exponentially distributed claim sizes, i.e. $F(u) = 1 - e^{-\frac{u}{\mu}}$, $u \geq 0$, $\mu > 0$. In this case,

$$h(x) = (1 - \rho) \frac{\alpha}{\mu(\alpha - t)} e^{-(1-\rho)\frac{\alpha x}{\mu(\alpha - t)}}$$

and

$$H(x) = 1 - e^{-(1-\rho)\frac{\alpha x}{\mu(\alpha - t)}}, \quad x \geq 0.$$

The first order differential equation (7) is given by

$$\begin{aligned} \frac{\partial G(u, y)}{\partial u} - \frac{n}{c(\alpha - t)} G(u, y) = \\ - \frac{n(1 - \rho)\alpha}{c\mu(\alpha - t)^2} e^{-(1-\rho)\frac{\alpha}{\alpha-t}\frac{y}{\mu}} \int_0^u G(v, y) e^{(1-\rho)\frac{\alpha}{\alpha-t}\frac{v}{\mu}} dv - \\ - \frac{n}{c(\alpha - t)} e^{-(1-\rho)\frac{\alpha}{\alpha-t}\frac{y}{\mu}} \left[1 - e^{-(1-\rho)\frac{\alpha}{\alpha-t}\frac{y}{\mu}} \right] \end{aligned}$$

Differentiating by u leads to the second order differential equation

$$\frac{\partial^2 G(u, y)}{\partial u^2} - \frac{n}{c(\alpha - t)} \left(1 - \frac{(1 - \rho)\alpha c}{n\mu} \right) \frac{\partial G(u, y)}{\partial u} = 0. \quad (10)$$

The initial condition (8) in the case of exponential distribution is

$$G(0, y) = \frac{n}{c(1 - \rho)\alpha} \left(1 - e^{-(1-\rho)\frac{\alpha}{\alpha-t}\frac{y}{\mu}} \right). \quad (11)$$

The equation (7) gives the second condition

$$\frac{\partial G(0, y)}{\partial u} = \frac{n\mu}{c(\alpha - t)} \left(\frac{n\mu}{c(1 - \rho)\alpha} - 1 \right) \left(1 - e^{-(1-\rho)\frac{\alpha}{\alpha-t}\frac{y}{\mu}} \right). \quad (12)$$

The solution of (10) with the initial conditions (11) and (12) is

$$G(u, y) = \frac{n\mu}{c(1 - \rho)\alpha} \left(1 - e^{-(1-\rho)\frac{\alpha}{\alpha-t}\frac{y}{\mu}} \right) e^{-\frac{n}{c(\alpha-t)} \left(\frac{c(1-\rho)\alpha}{n\mu} - 1 \right) u}.$$

The ruin probability in the exponential case is given by

$$\Psi(u) = \frac{n\mu}{c(1 - \rho)\alpha} e^{-\frac{n}{c(\alpha-t)} \left(\frac{c(1-\rho)\alpha}{n\mu} - 1 \right) u}.$$

References

1. Johnson, N.L., Kotz, S. and Kemp, A.W. (1992). *Univariate Discrete Distributions*, Wiley Series in Probability and Mathematical Statistics. 2nd edition.
2. Klugman S. A., Panjer H. and Willmot G. (1998) *Loss Models. From Data to Decisions*, John Wiley & Sons, Inc.
3. Minkova L.D. (2001). A Family of Compound Discrete Distributions, *Compt. Randue Bulg. Acad. Aci.*, 54(2), 11-14.
4. Minkova L.D. (2001). Inflated-parameter modification of the pure birth process, *Compt. Randue Bulg. Acad. Sci.* 54(11), 17-22.
5. Minkova L.D. (2002). A Generalization of the Classical Discrete Distributions, *Commun.Statist. - Theory and Methods*, 31(6), 871-888.

Families of Distributions Arising from Distributions of Record Statistics^{*}

S. M. T. K. Mirmostafae¹ and Jafar Ahmadi²

¹ Ferdowsi university of Mashhad, P. O. Box 91775-1159, Mashhad, Iran
(e-mail: ta_mi182@stu-mail.um.ac.ir)

² (e-mail: ahmadi-j@um.ac.ir)

Abstract. This paper proposed a new general family of continuous of distributions motivated by the distributions of record statistics. Its distributional properties including the distribution function, moments, symmetry and modality are studied. One special case, when F is the exponential distribution is considered and at the end two real data sets are used for fitting the suitability of our proposed model in the special case.

Keywords: Order statistics, Record statistics, Gamma distribution.

1 Introduction

Let $\{X_i, i \geq 1\}$ be a sequence of continuous random variables from the cumulative distribution function (cdf) $F(x)$ and probability density function (pdf) $f(x)$. Then the pdf of n -th upper record values, U_n , and n -th lower record values, L_n , are given by

$$f_{U_n}(x) = \frac{1}{\Gamma(n)} [-\log \bar{F}(x)]^{n-1} f(x) \quad -\infty < x < \infty \quad (1)$$

and

$$f_{L_n}(x) = \frac{1}{\Gamma(n)} [-\log F(x)]^{n-1} f(x) \quad -\infty < x < \infty \quad (2)$$

respectively, where $\Gamma(\cdot)$ is the complete gamma function. See Arnold *et al.* (1998) for more details about the theory and applications of record values. Several authors have considered the problems of generalized continuous probability distributions. The generalized gamma distribution, Pareto distribution and beta distribution have been studied by Amoroso (1925), Ljubo (1965) and McDonald (1984), respectively. Since then other authors have developed the previous results. Recently, Eugene *et al.* (2002) introduced a new family of distributions generated from the logit of the beta random variable:

$$g_F(t; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} [F(t)]^{\alpha-1} [1 - F(t)]^{\beta-1} f(t). \quad (3)$$

^{*} This paper is supported in part by the Iranian National Foundation of Elites.

They studied a special case of (3) when $F(x)$ is the cdf of the normal distribution. Jones (2004) proposed (3) as a family of distributions motivated by order statistics. He studied its distributional properties as well as potential for exciting statistical applications. Previous researches, (1) and (2) caused us to try to introduce a new family of distributions which arises from the distributions of record statistics. We propose the new general family of continuous distributions that generated by F as follows:

$$h_F(t; \alpha, \beta) = \frac{1}{\gamma(\alpha, \beta)} [-\log \bar{F}(t)]^{\alpha-1} [-\log F(t)]^{\beta-1} f(t), \quad (4)$$

where α and β are positive real constants and $\gamma(\alpha, \beta)$ is

$$\gamma(\alpha, \beta) = \int_0^\infty y^{\alpha-1} [-\log(1 - e^{-y})]^{\beta-1} e^{-y} dy, \quad (5)$$

we call it extended gamma function. It is clear that $\gamma(\alpha, 1) = \Gamma(\alpha)$ and $\gamma(1, \beta) = \Gamma(\beta)$. From (1), (2) and (4) it is obvious that if $\beta = 1$ and $\alpha \in N$, as a natural number, then the probability distribution in (4) takes the same form as the α -th upper record values. Also for $\alpha = 1$ and $\beta \in N$ the pdf in (4) is the pdf of lower record values coming from cdf F and pdf f . Thus in the case of $\alpha = 1$, one example of family (4) is the gamma distribution itself which arises immediately if F is taken to be the exponential distribution. Also, $h_F(x; 1, 1) = f(x)$. The main reason for extending a distribution as in (4) is that the form in (4) provides more flexibility in modelling observed data.

We study some distributional properties of the introduced family in Section 2. These properties include the distribution function, moments, symmetry, modality and estimation of α and β . One special case, when F is the exponential distribution is considered in Section 3. At the end two real data sets are used for fitting the suitability of our proposed model in the special case.

2 Distribution Properties

In this section we intend to study the general properties of members of family (4). First of all we present some properties of $\gamma(\alpha, \beta)$.

2.1 Some properties of $\gamma(\alpha, \beta)$

Lemma 1. *For any positive real values of α and β we have:*

- $[\gamma(\alpha, \beta)]^2 \leq \Gamma(2\alpha - 1)\Gamma(2\beta - 1)$ and the equality holds iff $\alpha = \beta = 1$.
- $\gamma(\alpha, \beta) = \gamma(\beta, \alpha)$.

The proof of the above lemma is simple and therefore is omitted. From (5) it is obvious that $\gamma(\alpha, 1) = \Gamma(\alpha)$ and in the following lemma, we obtain the exact expression for $\gamma(\alpha, 2)$ which will be used in this paper.

Lemma 2. For $\beta = 2$ and $\alpha > 0$, we have

$$\gamma(\alpha, 2) = \Gamma(\alpha) \sum_{j=1}^{\infty} \frac{1}{j(j+1)^\alpha}.$$

Proof. The result immediately follows by notifying that

$$\log(1 - e^{-y}) = - \sum_{i=1}^{\infty} \frac{e^{-iy}}{i}.$$

2.2 Distribution Function

From (4) and putting $y = -\log \bar{F}(x)$ we have

$$H_F(x; \alpha, \beta) = \frac{\gamma(\alpha, \beta, -\log \bar{F}(x))}{\gamma(\alpha, \beta)},$$

where $\gamma(\alpha, \beta, t)$ is the incomplete extended gamma function, i.e.,

$$\gamma(\alpha, \beta, t) = \int_0^t y^{\alpha-1} [-\log(1 - e^{-y})]^{\beta-1} e^{-y} dy.$$

Now, we calculate $H_F(x, \alpha, \beta)$ for some special cases.

- Suppose $\beta - 1$ is a natural number ($\beta - 1 \in N$), then

$$\begin{aligned} H_F(x; \alpha, \beta) &= \frac{\gamma(\alpha, \beta, -\log \bar{F}(x))}{\gamma(\alpha, \beta)} \\ &= \gamma^{-1}(\alpha, \beta) \int_0^{-\log \bar{F}(x)} y^{\alpha-1} [-\log(1 - e^{-y})]^{\beta-1} e^{-y} dy \\ &= \gamma^{-1}(\alpha, \beta) \int_0^{-\log \bar{F}(x)} y^{\alpha-1} \left(\sum_{i=1}^{\infty} \frac{e^{-iy}}{i} \right)^{\beta-1} e^{-y} dy \\ &= \gamma^{-1}(\alpha, \beta) \int_0^{-\log \bar{F}(x)} y^{\alpha-1} \sum_{j=1}^{\infty} C_j(\beta - 1) e^{-(j+1)y} dy \\ &= \gamma^{-1}(\alpha, \beta) \sum_{j=1}^{\infty} C_j(\beta - 1) \Gamma(\alpha, j + 1, -\log \bar{F}(x)), \end{aligned} \tag{6}$$

where $C_j(n)$ in (6) is the coefficient of w^j in the expansion of $\left(\sum_{i=1}^{\infty} \frac{w^i}{i}\right)^n$ and $\Gamma(\alpha, \beta, t)$ is the incomplete gamma function. Notice that the coefficients $C_j(n)$ can be generated in a recursive manner as follows: $C_j(1) = 1/j$ for $j = 1, 2, \dots$ and

$$C_j(n) = \sum_{k=n-1}^{j-1} C_k(n-1)/(j-k),$$

see Arnold *et al.* (1998), pp. 70–71.

- Suppose that α is a natural number in (6), by using the following identity

$$\int_0^t \frac{\lambda^n}{\Gamma(n)} x^{n-1} e^{-\lambda x} dx = \sum_{k=n}^{\infty} \frac{(\lambda t)^k}{k!} e^{-\lambda t}.$$

We can rewrite (6) as follows

$$H_F(x; \alpha, \beta) = \gamma^{-1}(\alpha, \beta) \Gamma(\alpha) \sum_{j=1}^{\infty} \sum_{i=\alpha}^{\infty} \frac{C_j(\beta-1)}{i! j^{\alpha-i}} [\bar{F}(x)]^j [-\log \bar{F}(x)]^i.$$

For $\beta = 2$ and $\alpha \in N$, $H_F(x, \alpha, 2)$ simplifies as:

$$H_F(x; \alpha, 2) = \left(\sum_{k=1}^{\infty} \frac{1}{k(k+1)^\alpha} \right)^{-1} \sum_{j=1}^{\infty} \sum_{i=\alpha}^{\infty} \frac{[\bar{F}(x)]^{j+1} [-\log \bar{F}(x)]^i}{j(j+1)^{\alpha-i} i!}.$$

2.3 Moments

A sufficient condition for existence of the moments of the family in (4) is provided by the following lemma. We say that the k th moment of X exists if $E(|X|^k) < \infty$.

Lemma 3. *Let U have a distribution with pdf (4) and X have a distribution with cdf $F(x)$. If $E(|X|^{k+\delta}) < \infty$, k is any non-negative integer and $\delta > 0$, then $E(|U|^k) < \infty$, for all $\alpha > 0$, $\beta > 0$.*

Proof. Suppose that $p > 1$ and $\frac{1}{p} + \frac{1}{q} = 1$, then from (4) and taking $y = -\log \bar{F}(u)$ we have

$$\begin{aligned} E(|U|^k) &= \gamma^{-1}(\alpha, \beta) \int_{-\infty}^{\infty} |u|^k [-\log \bar{F}(u)]^{\alpha-1} [-\log F(u)]^{\beta-1} f(u) du \\ &= \gamma^{-1}(\alpha, \beta) \int_0^{\infty} |F^{-1}(1 - e^{-y})|^k [-\log(1 - e^{-y})]^{\beta-1} e^{-y} y^{\alpha-1} dy \\ &\leq \gamma^{-1}(\alpha, \beta) \left\{ \int_0^{\infty} \{y^{\alpha-1} [-\log(1 - e^{-y})]^{\beta-1}\}^q e^{-y} dy \right\}^{\frac{1}{q}} \\ &\quad \times \left\{ \int_0^{\infty} |F^{-1}(1 - e^{-y})|^{kp} e^{-y} dy \right\}^{\frac{1}{p}} \tag{7} \\ &= \frac{[\gamma(\alpha q - q + 1, \beta q - q + 1)]^{\frac{1}{q}}}{\gamma(\alpha, \beta)} \left\{ \int_{-\infty}^{\infty} |x|^{k+\delta} f(x) dx \right\}^{\frac{1}{p}}, \text{ where } k + \delta = kp \\ &= \frac{[\gamma(\alpha q - q + 1, \beta q - q + 1)]^{\frac{1}{q}}}{\gamma(\alpha, \beta)} [E(|X|^{k+\delta})]^{1/p}, \end{aligned}$$

The inequality in (7) Holder’s inequality.

So by Lemma 3 the existence of the moments of F will guarantee existence of the moments of lower order of the corresponding generated family in (4).

2.4 Symmetry and Modality

In this section we seek to provide a sufficient condition for preserving symmetry and unimodality properties by h_F .

Lemma 4. *Let F be symmetric about zero, then h_F remains symmetric whenever $\alpha = \beta$.*

The proof of the above lemma is simple and therefore is omitted. It is not difficult to show that $\gamma(\alpha, \alpha)$ in (5) is a increasing function with respect to α . So whenever $\alpha = \beta$, h_F remains symmetric but with tails getting lighter as α increases and heavier as α decreases. If $\alpha \neq \beta$, skewness is introduced, the amount of skewness depends on the difference between α and β , and its sign on the sign of $\beta - \alpha$.

Lemma 5. *Let F be symmetric and unimodal, then h_F is also unimodal, if $\alpha = \beta$.*

Proof. Let M be the mode of F , then by assumptions $F(M) = \frac{1}{2}$. From (4) for $\alpha = \beta$,

$$h_F(x; \alpha, \alpha) = \gamma^{-1}(\alpha, \alpha) [\log \bar{F}(x) \log F(x)]^{\alpha-1} f(x). \tag{8}$$

So, it is enough to show that the expression in the bracket on the right hand side of (8) gets its only maximum at $x = M$. To this end, let $g(x) = \log \bar{F}(x) \log F(x)$, then

$$g'(x) = \frac{f(x)}{F(x)\bar{F}(x)} [\bar{F}(x) \log \bar{F}(x) - F(x) \log F(x)],$$

and

$$g''(x) = -S(x)\{f(x)[2 + \log(F(x)\bar{F}(x))]\} + R(x)S'(x),$$

where $S(x) = \frac{f(x)}{F(x)\bar{F}(x)}$ and $R(x) = \bar{F}(x) \log \bar{F}(x) - F(x) \log F(x)$. We have

$$\lim_{|x| \rightarrow \infty} g'(x) = \lim_{|x| \rightarrow \infty} R(x) = 0, \tag{9}$$

and as $S(x) > 0$, $g'(x) = 0$ whenever $R(x) = 0$. Clearly $R(M) = 0$, and $g''(M) < 0$, so by (9) it is enough to show that $R'(x)$ has only two finite roots. We have

$$R'(x) = -f(x)[2 + \log(F(x)\bar{F}(x))]. \tag{10}$$

The right hand side of (10) equals to zero whenever $x = x_1 = F^{-1}\left(\frac{2e^{-1}}{\sqrt{e^2-4}+e}\right)$ or $x = x_2 = F^{-1}\left(\frac{2e^{-1}}{\sqrt{e^2-4}-e}\right)$. It is clear that $M \in (x_1, x_2)$ and the proof is complete.

In Lemma 4 and 5, we show that the properties of symmetry and unimodality preserve by h_F when $\alpha = \beta$. The parameters α and β are shape parameters, which determines the skewness of the distribution. When F is unimodal, h_F is skewed to the right when $\beta > \alpha$, the degree of right skewness increases as β increases. Also h_F is skewed to the left when $\beta < \alpha$, the degree of left skewness increases as β decreases.

2.5 Estimation of α and β

Let F be free of parameters and suppose X_1, X_2, \dots, X_n constitute a random sample of size n from (4), then the likelihood function is given by:

$$L(\alpha, \beta) = [\gamma(\alpha, \beta)]^{-n} \prod_{i=1}^n f(x_i) \exp\{(\alpha-1) \log(-\log \bar{F}(x_i)) + (\beta-1) \log(-\log F(x_i))\}.$$

Then, clearly $W_1(\mathbf{X}) = \sum_{i=1}^n \log(-\log \bar{F}(X_i))$ and $W_2(\mathbf{X}) = \sum_{i=1}^n \log(-\log F(X_i))$ are complete sufficient statistics for α and β , respectively. The maximum likelihood estimator (MLE) of α and β can be obtained by solving the following two equations:

$$W_1(\mathbf{x}) - n \frac{\partial \gamma(\alpha, \beta)}{\partial \alpha} / \gamma(\alpha, \beta) = 0, \quad \text{and} \quad W_2(\mathbf{x}) - n \frac{\partial \gamma(\alpha, \beta)}{\partial \beta} / \gamma(\alpha, \beta) = 0.$$

The Fisher Information matrix for (α, β) is given by

$$I(\alpha, \beta) = \begin{bmatrix} \frac{\partial^2 \gamma(\alpha, \beta)}{\partial \alpha^2} & \frac{\partial^2 \gamma(\alpha, \beta)}{\partial \alpha \partial \beta} \\ \frac{\partial^2 \gamma(\alpha, \beta)}{\partial \alpha \partial \beta} & \frac{\partial^2 \gamma(\alpha, \beta)}{\partial \beta^2} \end{bmatrix}.$$

If one of the parameters is known, say β , then it is obvious that $I_X(\alpha) = \frac{\partial^2 \gamma(\alpha, \beta)}{\partial \alpha^2}$ and by MLE properties

$$\sqrt{n}(\hat{\alpha} - \alpha) \longrightarrow N(0, I_X^{-1}(\alpha)).$$

3 Special Case (Extended gamma distribution)

As pointed in section 1, recently several new distributions were introduced in the literature. Here, we consider a special case of (4), $\beta = 2$ and suppose X has exponential distribution, i.e. $\bar{F}(x) = \exp(-\lambda x)$ $\lambda > 0$, in our model. Then from (4) we find

$$h_F(x; \alpha, 2) = \frac{\lambda^\alpha}{\eta(\alpha)} x^{\alpha-1} e^{-\lambda x} [-\log(1 - e^{-\lambda x})], \quad x > 0, \quad (11)$$

where from Lemma 2, $\eta(\alpha) = \gamma(\alpha, 2) = \Gamma(\alpha) \sum_{j=1}^{\infty} \frac{1}{j(j+1)^\alpha}$. Here α is the shape parameter and λ is the scale parameter. When $\alpha = 1$, the model (11) follows the distribution of the second lower record values from the exponential distribution. We say that random variable X has extended gamma distribution (EG) and denote $X \sim EG(\alpha, \lambda)$, if its pdf is as in (11). Then we have

$$E(X^k) = \frac{\eta(\alpha + k)}{\eta(\alpha) \lambda^k}, \quad k \geq 0.$$

So the moment estimators of λ and α can be obtained as $\hat{\lambda} = \frac{\eta(\hat{\alpha}+1)}{\eta(\hat{\alpha})\bar{X}}$, where $\hat{\alpha}$ satisfies the following identity

$$\frac{\eta(\hat{\alpha} + 2)\eta(\hat{\alpha})}{[\eta(\hat{\alpha} + 1)]^2} = \bar{X}^2 \bar{X}^2.$$

Let X_1, X_2, \dots, X_n be a random sample of size n from EG, then the log likelihood function can be written as:

$$\begin{aligned} L(\alpha, \lambda) = & -n \log \eta(\alpha) + (\alpha - 1) \sum_{i=1}^n \log x_i + n \alpha \log \lambda \\ & + \sum_{i=1}^n \log (-\log(1 - e^{-\lambda x_i})) - \lambda \sum_{i=1}^n x_i. \end{aligned}$$

On taking partial derivatives of the log likelihood with respect to α and λ respectively and equating the derivatives to zero we get

$$\begin{aligned} \frac{\partial L}{\partial \alpha} = & -n \frac{\partial \eta(\alpha)}{\partial \alpha} + \sum_{i=1}^n \log x_i + n \log \lambda = 0, \\ \frac{\partial L}{\partial \lambda} = & \frac{n \alpha}{\lambda} + \sum_{i=1}^n \frac{x_i e^{-\lambda x_i}}{(1 - e^{-\lambda x_i}) \log(1 - e^{-\lambda x_i})} - \sum_{i=1}^n x_i = 0. \end{aligned}$$

Therefore, we can obtain the MLE's of α and λ by solving the above non-linear normal equations. From the second equation $\hat{\alpha}$ can be obtained as a function of λ as follows:

$$\hat{\alpha}(\lambda) = \frac{\lambda}{n} \left(\sum_{i=1}^n x_i - \sum_{i=1}^n \frac{x_i e^{-\lambda x_i}}{(1 - e^{-\lambda x_i}) \log(1 - e^{-\lambda x_i})} \right). \tag{12}$$

Let $\varphi(\alpha) = \frac{\partial}{\partial \alpha} \eta(\alpha)$. So if both of the parameters are unknown, first the MLE λ , say $\hat{\lambda}$, can be obtained by maximizing directly

$$\begin{aligned} g(\lambda) = L(\hat{\alpha}(\lambda), \lambda) = & -n \varphi \left(\frac{\lambda}{n} \left(\sum_{i=1}^n x_i - \sum_{i=1}^n \frac{x_i e^{-\lambda x_i}}{(1 - e^{-\lambda x_i}) \log(1 - e^{-\lambda x_i})} \right) \right) \\ & + \sum_{i=1}^n \log x_i + n \log \lambda. \end{aligned}$$

with respect to λ . Once $\hat{\lambda}$ is obtained, $\hat{\alpha}$ can be obtained from (12).

3.1 Data Analysis

In order to fit EG to data, we used two real data sets represent the failure times of the air conditioning systems of two different air planes (see Bain

and Engelhart, 1991). Gupta and Kundu (2003) fitted both the gamma distribution and exponentiated exponential (EE) distribution to these data.

Data 1: Plane 7912: 1, 3, 5, 7, 11, 11, 11, 12, 14, 14, 14, 16, 16, 20, 21, 23, 42, 47, 52, 62, 71, 71, 87, 90, 95, 120, 120, 225, 246, 261.

Data 2: Plane 7911: 33, 47, 55, 56, 104, 176, 182, 220, 239, 246, 320.

We fit gamma, EE and EG distribution functions to these data. We also estimate the unknown parameters in all these cases by maximum likelihood method. Moreover we present the χ^2 statistics for these three cases. The results are summarized in Table 1. From Table 1 it is observed that, by

Data set	Distribution	$\hat{\lambda}$	$\hat{\alpha}$	χ^2
1	Gamma	0.0136	0.8134	3.302
	EE	0.0145	0.8130	3.383
	EG	0.0066	1.0766	3.181
2	Gamma	0.014	2.1457	0.9929
	EE	0.104	2.2427	1.0917
	EG	0.007	2.5136	0.9928

Table 1. The goodness of fit of gamma, EE and EG distributions to data sets 1 and 2.

empirical evidence, in both cases the EG distribution is fitted better than that gamma and EE distributions. It may be noted that in the second data set the χ^2 statistic for EG distribution is very close to that for gamma distribution. Notice that the results of this section not guarantee that the EG will always better than EE or gamma distributions, but at least it can be said that in some cases, it is better. One or two examples do not tell us much more.

References

1. Amoroso, L. "Ricerche intorno alla curva dei redditi". *Annali de Mathematica Seies* 4, 2:123–159 (1925).
2. Arnold, B. C., Balakrishnan, N. and Nagaraja, H. N., *Records*, Wiley, New York (1998).
3. Bain, L. J. and Engelhardt, M. *Statistical Analysis of Reliability and life testing Models*. Second, Edition: Marcel and Dekker, New York (1991).
4. Eugene, N., Lee, C. and Famoye, F. "Beta-normal distribution and its applications". *Comm. Statist. Theory Methods*, 31:497–512 (2002).
5. Gupta, R. D. and Kundu, D. "Closeness of gamma and generalized exponential distribution". *Comm. Statist. Theory Methods*, 32:705–721 (2003).
6. Jones, M. C. "Families of distributions arising from distributions of order statistics". *Test*, 13:1–43 (2004).
7. Ljubo, M. "Curves and concentration indices for certain generalized Pareto distributions". *Statisti cal Review*, 15:257–260 (1965).
8. McDonald, J.B. "Some generalized functions for the size distribution of income". *Econometrica*, 52:647–663 (1984).

Extinction Probability in the Class of Two-Sex Branching Models with Offspring and Mating Depending on the Number of Couples

Manuel Molina¹, Yongsheng Xing², and Shixia Ma³

¹ Department of Mathematics, University of Extremadura, Badajoz, Spain
(e-mail: mmolina@unex.es)

² College of Mathematics, Shandong Institute of technology, Yantai, China
(e-mail: xingys@nankai.edu.cn)

³ School of Sciences, Hebei University of Technology, Tianjin, China
(e-mail: mashixia1@163.com)

Abstract. This paper deals with stochastic modeling through branching models. It is our purpose to model the probabilistic evolution of populations where females and males coexist and form couples. In particular, the class of two-sex branching models with offspring and mating depending on the number of couples in the population is considered. This class has practical implications, especially in population dynamics. For such a class of models, by considering different approaches, we provide some necessary and sufficient conditions for the almost sure extinction of the process.

Keywords: Branching models, Two-sex models, Extinction probability.

1 Introduction

With the purpose to model the probabilistic evolution of populations where females and males coexist and form couples (female-male) several classes of discrete time branching models have been investigated, including the bisexual Galton-Watson model (see Alsmeyer and Rösler (1996) [1], (2002) [2], Bruss (1984) [3], Daley (1968) [4], Daley *et al.* (1986) [5]), two-sex models with immigration (see González *et al.* (2000) [6], (2001) [7], Ma and Xing (2006) [10]), in varying environments (see Molina *et al.* (2003) [13]), in random environments (see Ma (2006) [8], Ma and Molina (2009) [9]), with population-size depending mating (see Molina *et al.* (2002) [12], (2004) [14], (2006) [15], Xing (2005) [18]), or with a control function (see Molina *et al.* [16]). Recently, it has been introduced, see Molina *et al.* (2008) [11], the class of two-sex branching models with offspring and mating depending on the number of couples in the population. Several relationships among the probability generating functions involved in the stochastic model have been determined and some limiting results derived. The aim of this paper is to continue the research about such a class of two-sex models, investigating necessary and sufficient conditions for its almost sure extinction.

The paper is organized as follow. In Section 2, the two-sex process is described formally and interpreted intuitively. Section 3 is devoted to determining some results concerning the extinction probability of the model. Finally, the proofs are included in Section 4.

2 The two-sex model

Let us consider the two-sex branching model $\{(F_n, M_n)\}_{n \geq 1}$ defined in the form:

$$(F_n, M_n) = \sum_{i=1}^{Z_{n-1}} (f_{n,i}(Z_{n-1}), m_{n,i}(Z_{n-1})), \quad Z_n = L_{Z_{n-1}}(F_n, M_n), \quad n = 1, 2, \dots \quad (1)$$

where the empty sum is considered to be $(0, 0)$. The random vector (F_n, M_n) represents the number of females and males in the n th generation. These females and males form Z_n couples. A couple consists of a female and a male from the same generation who came with the purpose of generating offspring. It is assumed that initially there are $N_0 \geq 1$ couples in the population, i.e., $Z_0 = N_0$. Let us denote by \mathbb{Z}^+ and \mathbb{R}^+ , respectively, the non-negative integer and real numbers. Given that, in the $(n-1)$ th generation there are N couples, namely $Z_{n-1} = N$, then:

- (a) L_N is the function which governs the mating between females and males. It is a non-negative real function, defined on $\mathbb{R}^+ \times \mathbb{R}^+$, assumed to be non-decreasing in each argument, integer-valued on the integers, and such that, for $x, y \in \mathbb{R}^+$, $L_N(x, 0) = L_N(0, y) = 0$.
- (b) $\{(f_{n,i}(N), m_{n,i}(N)), i = 1, \dots, N\}$ are independent and identically distributed non-negative, integer-valued random vectors. Intuitively, the random vector $(f_{n,i}(N), m_{n,i}(N))$ represents the number of females and males descending from the i th couple of the $(n-1)$ th generation. Its probability law will be referred as the offspring probability distribution when there are N progenitor couples in the population. Clearly, $P(f_{1,1}(0) = 0, m_{1,1}(0) = 0) = 1$.

Note that $\{(F_n, M_n)\}_{n \geq 1}$ may be interpreted as a stochastic model developing in an environment which changes in time according to the number of couples in the population. In each generation, both the offspring probability distribution and the mating function are affected by the number of couples in the previous generation. In addition to its theoretical interest, this class of two-sex models also has practical implications, especially in population dynamics. In facts, by environmental, social, or other factors, the offspring and the mating between females and males may be affected by the number of couples in the population. Indeed, the motivation behind this class of processes has been the interest in developing models to describe such behaviors.

The class of models given in (1) includes, as particular cases, the two-sex models introduced in Daley (1968) [5], Molina *et al.* (2002) [12], and Xing and Wang (2005) [18].

In order to establish some results about its extinction probability, we shall consider the following requirements about the mating functions and the offspring probability distributions:

(a1): $\{L_N\}_{N \geq 0}$ is such that L_N is a superadditive function, namely,

$$L_N(x_1 + x_2, y_1 + y_2) \geq L_N(x_1, y_1) + L_N(x_2, y_2), \quad x_i, y_i \in \mathbb{R}^+, \quad i = 1, 2.$$

(a2) $\{L_N(x, y)\}_{N \geq 0}$, where $x, y \in \mathbb{R}^+$ are fixed, is a non-decreasing sequence.

(a3) $f_{1,1}(N) \preceq^1 f_{1,1}(N + 1)$; $m_{1,1}(N) \preceq m_{1,1}(N + 1)$, $N \in \mathbb{Z}^+$.

Remark 1. Assumption (a1) expresses the intuitive notion that $x_1 + x_2$ females and $y_1 + y_2$ males coexisting together will form a number of couples that is at least as great as the total number of couples formed by x_1 females and y_1 males, and x_2 females and y_2 males, living separately. Most of mating functions considered in two-sex branching model theory are super-additive. Assumption (a2) represents the usual behavior in many biological populations in which the mating is promoted as the number of couples grows. According to (a3), the variables $f_{1,1}(N)$ and $m_{1,1}(N)$ take large values with a lower probability than $f_{1,1}(N + 1)$ and $m_{1,1}(N + 1)$ do, respectively. This expresses the intuitive fact that when the number of couples in the population grows then the corresponding numbers of originated females and males take large values with higher probabilities.

Throughout this work, we will assume the classical duality extinction-explosion in branching model theory, namely, for $N \geq 1$,

$$P(\lim_{n \nearrow \infty} Z_n = 0 \mid Z_0 = N) + P(\lim_{n \nearrow \infty} Z_n = \infty \mid Z_0 = N) = 1. \quad (2)$$

Under this framework, some sufficient conditions for the non-extinction of $\{(F_n, M_n)\}_{n \geq 1}$ have been determined in [11] where some general settings which guarantee (2) holds have been established. Also, it has been proved that $R := \lim_{N \rightarrow \infty} R_N$ exists where

$$R_N = N^{-1}E[Z_n \mid Z_{n-1} = N], \quad N = 1, 2, \dots$$

Next, we continue the research about the extinction probability concerning such a class of two-sex branching models.

¹ Given the random variables X and Y , we say that X is stochastically smaller than Y , written $X \preceq Y$, if $P(X > t) \leq P(Y > t)$, $t \in \mathbb{R}$.

3 Extinction probability

In this section we shall derive some necessary and sufficient conditions concerning the extinction probability of the class of models presented in (1). To this end, we shall use two different approaches. First, by considering the asymptotic growth rate R (Proposition 1) and then, by using the comparison with a simpler two-sex model (Proposition 2).

Remark 2. If for some $n \geq 1$, $Z_n = 0$ then, from (1), one deduces that $Z_{n+m} = 0$ and $(F_{n+m}, M_{n+m}) = (0, 0)$, $m \geq 1$. Hence the two-sex model does not survive.

Definition 1. For every $N \geq 1$, let

$$Q_N = P(\lim_{n \nearrow \infty} Z_n = 0 \mid Z_0 = N)$$

be the extinction probability when initially there are N couples in the population.

Proposition 1. Assume (a1), (a2) and (a3).

(i) If $R \leq 1$ then $Q_N = 1$ for $N \geq 1$.

(ii) If $R > 1$ then there exists $K_0 \geq 1$ such that $Q_N < 1$ for $N \geq K_0$.

Remark 3. In the following result, by using a methodology based in the stochastic comparison with a two-sex model with only mating depending on the number of couples in the population, necessary and sufficient condition for the almost sure extinction are also determined. First, we shall introduce the following modification in requirement (a3):

(a4): For $N \in \mathbb{Z}^+$, $f_{1,1}(N) \preceq f_{1,1}(N+1)$, $m_{1,1}(N) \preceq m_{1,1}(N+1)$ and there exist random variables $f_{1,1}$ and $m_{1,1}$ such that $\lim_{N \nearrow \infty} f_{1,1}(N) = f_{1,1}$ and $\lim_{N \nearrow \infty} m_{1,1}(N) = m_{1,1}$ almost surely.

Taking into account (a4), one deduces,

$$f_{1,1}(N) \preceq f_{1,1}, \quad m_{1,1}(N) \preceq m_{1,1}, \quad N \in \mathbb{Z}^+$$

Let $(\mu_f(N), \mu_m(N))$ and (μ_f, μ_m) be the mean vectors of $(f_{1,1}(N), m_{1,1}(N))$ and $(f_{1,1}, m_{1,1})$, respectively, both assumed to be finite. Again, by (a4), one derives that $\{\mu_f(N)\}_{N \geq 0}$ and $\{\mu_m(N)\}_{N \geq 0}$ are non-decreasing sequences. Hence, by the monotone convergence theorem,

$$\lim_{N \nearrow \infty} \mu_f(N) = \mu_f, \quad \lim_{N \nearrow \infty} \mu_m(N) = \mu_m$$

Let $\{(F_n^*, M_n^*)\}_{n \geq 1}$ be the two-sex model initiated with $Z_0^* = N_0$:

$$(F_n^*, M_n^*) = \sum_{i=1}^{Z_{n-1}^*} (f_{n,i}, m_{n,i}), \quad Z_n^* = L_{Z_{n-1}^*}(F_n^*, M_n^*), \quad n = 1, 2, \dots \quad (3)$$

where $\{(f_{n,i}, m_{n,i})\}_{n,i \geq 1}$ is a sequence of independent and identically distributed random vectors with the same probability law of $(f_{1,1}, m_{1,1})$. Model (3) was introduced in Molina *et al.* (2002) [12] where it was established that $R^* := \lim_{k \nearrow \infty} R_k^* = \sup_{k > 0} R_k^*$ exists, $R_k^* = k^{-1}E[Z_n^* | Z_{n-1}^* = k]$, $k \geq 1$ and, moreover, $R^* \leq 1$ if and only if

$$Q_N^* := P(\lim_{n \nearrow \infty} Z_n^* = 0 | Z_0^* = N) = 1, \quad N \geq 1$$

Proposition 2. *Assume (a1), (a2) and (a4).*

- (i) *If $R^* \leq 1$ then $Q_N = 1$ for $N \geq 1$.*
- (ii) *If $R^* > 1$ then there exists $K_0 \geq 1$ such that $Q_N < 1$ for $N \geq K_0$.*

4 Proofs

4.1 Proof of Proposition 1

By using (a1), (a2), and (a3), it is deduced, see Molina et al.(2008)[11], that $R = \sup_{N > 0} R_N$.

- (i) If $R \leq 1$ then $\{E[Z_n]\}_{n \geq 0}$ is a non-increasing sequence. In fact,

$$E[Z_{n+1}] = E[E[Z_{n+1} | Z_n]] = E[Z_n R_{Z_n}] \leq E[Z_n R] \leq E[Z_n], \quad n \in \mathbb{Z}^+.$$

Hence,

$$P(\lim_{n \nearrow \infty} Z_n = \infty | Z_0 = N) = 0, \quad N \geq 1$$

and, by (2), one has that $q_N = 1$, $N \geq 1$.

- (ii) Assume $R > 1$. Since $R = \lim_{N \nearrow \infty} R_N$, there exists $K > 0$ such that for $N \geq K$, $R_N > 1$. Let us consider the auxiliary process: $\{(F'_n, M'_n)\}_{n \geq 1}$,

$$(F'_n, M'_n) = (F_n, M_n)I_{\{Z'_{n-1} \leq K\}} + \sum_{i=1}^{Z'_{n-1}} (f_{n,i}(K), m_{n,i}(K))I_{\{Z'_{n-1} > K\}},$$

$$Z'_n = Z_n I_{\{Z'_{n-1} \leq K\}} + L_K(F'_n, M'_n)I_{\{Z'_{n-1} > K\}}, \quad n = 1, 2, \dots$$

where $Z'_0 = N_0$ and I_A denotes the indicator function of the set A . It is verified that $Z'_n \preceq Z_n$, $n \in \mathbb{Z}^+$. Hence, taking into account Müller and Stoyan (2002), p. 3 [17], it is derived that, for $N \geq 1$,

$$P(\lim_{n \nearrow \infty} Z_n = \infty \mid Z_0 = N) \geq P(\lim_{n \nearrow \infty} Z'_n = \infty \mid Z'_0 = N). \quad (4)$$

Let $\{(F_n^{(K)}, M_n^{(K)})\}_{n \geq 1}$ be the bisexual Galton-Watson process initiated with $Z_0^{(K)} = N_0$ couples and defined, for $n=1,2,\dots$, in the form:

$$(F_n^{(K)}, M_n^{(K)}) = \sum_{i=1}^{Z_{n-1}^{(K)}} (f_{n,i}(K), m_{n,i}(K)), \quad Z_n^{(K)} = L_K(F_n^{(K)}, M_n^{(K)})$$

By Daley et al. (1986) [5], one deduces that

$$R^{(K)} := \lim_{N \nearrow \infty} R_N^{(K)} = \sup_{N > 0} R_N^{(K)}.$$

Clearly $R^{(K)} \geq R_K^{(K)}$. Now,

$$R_K^{(K)} = K^{-1} E [Z_n^{(K)} \mid Z_{n-1}^{(K)} = K] = K^{-1} E [Z_n \mid Z_{n-1} = K] = R_K > 1.$$

Thus, by bisexual Galton-Watson process theory, one deduces the existence of $K^* \in \mathbb{Z}^+$ such that, for $N \geq K^*$

$$P\left(\lim_{n \nearrow \infty} Z_n^{(K)} = \infty \mid Z_0^{(K)} = N\right) > 0.$$

Let $K_0 := \max\{K, K^*\}$. Then,

$$P\left(\lim_{n \nearrow \infty} Z_n^{(K)} = \infty \mid Z_0^{(K)} = K_0\right) > 0$$

and using the fact that $\{Z_n^{(K)}\}_{n \geq 0}$ is a homogeneous Markov chain,

$$P\left(\lim_{n \nearrow \infty} Z_n^{(K)} = \infty, Z_n^{(K)} \geq K, n > 0 \mid Z_0^{(k)} = K_0\right) > 0. \quad (5)$$

Hence, by comparing $\{Z'_n\}_{n \geq 0}$ and $\{Z_n^{(K)}\}_{n \geq 0}$ and by (5),

$$P(\lim_{n \nearrow \infty} Z'_n = \infty \mid Z'_0 = K_0) > 0. \quad (6)$$

Finally, from (4) and (6),

$$P(\lim_{n \nearrow \infty} Z_n = \infty \mid Z_0 = N) > 0, \quad N \geq K_0.$$

By (2), one obtains that $q_N < 1$ for $N \geq K_0$.

4.2 Proof of Proposition 2

From Proposition 1, it is sufficient to prove that $R^* = R$. By (a1), (a2) and (a4), the existence of R is assumed.

For each $N \in \mathbb{Z}^+$, let $\{(F_n^{(N)}, M_n^{(N)})\}_{n \geq 1}$ be the process, initiated with $Z_0^{(N)} = N_0$, and defined, for $n \geq 1$:

$$(F_n^{(N)}, M_n^{(N)}) = \sum_{i=1}^{Z_{n-1}^{(N)}} (f_{n,i}(N), m_{n,i}(N)), \quad Z_n^{(N)} = L_{Z_{n-1}^{(N)}}(F_n^{(N)}, M_n^{(N)}) \quad (7)$$

Process (7) is again a two-sex model with only mating depending on the number of couples, being the offspring distribution the probability law of $(f_{1,1}(N), m_{1,1}(N))$. Hence, for $N \in \mathbb{Z}^+$, there exists $R^{(N)} := \lim_{k \nearrow \infty} R_k^{(N)}$ and

$$R^{(N)} = \sup_{k > 0} R_k^{(N)}, \quad R_k^{(N)} = k^{-1} E[Z_n^{(N)} \mid Z_{n-1}^{(N)} = k], \quad k = 1, 2, \dots$$

Now, from (a4), taking into account stochastic order properties,

$$\sum_{i=1}^N f_{n,i}(N) \preceq \sum_{i=1}^N f_{n,i}, \quad \sum_{i=1}^N m_{n,i}(N) \preceq \sum_{i=1}^N m_{n,i}$$

and

$$E \left[L_N \left(\sum_{i=1}^N f_{n,i}(N), \sum_{i=1}^N m_{n,i}(N) \right) \right] \leq E \left[L_N \left(\sum_{i=1}^N f_{n,i}, \sum_{i=1}^N m_{n,i} \right) \right]$$

Therefore

$$R = \limsup_{N \nearrow \infty} R_N \leq \limsup_{N \nearrow \infty} R_N^* = R^*.$$

On the other hand, given $j \geq 1$ fixed, one derives for $N \geq j$,

$$E \left[L_N \left(\sum_{i=1}^N f_{n,i}(N), \sum_{i=1}^N m_{n,i}(N) \right) \right] \geq E \left[L_N \left(\sum_{i=1}^N f_{n,i}(j), \sum_{i=1}^N m_{n,i}(j) \right) \right].$$

Thus

$$R = \liminf_{N \nearrow \infty} R_N \geq \liminf_{N \nearrow \infty} R_N^{(j)} = R^{(j)}$$

Taking limit as $j \nearrow \infty$, one derives that $R \geq \lim_{j \nearrow \infty} R^{(j)}$. Finally, it is matter of straightforward calculation to deduce that $\lim_{j \nearrow \infty} R^{(j)} = R^*$, and consequently the proof is completed.

Acknowledgement

This research has been supported by the Ministerio de Ciencia e Innovación of Spain, grant MTM2009-13248, and by the Natural Sciences Foundation of China, grant 10971048.

References

1. Alsmeyer, G, and Rösler, U. “The bisexual Galton-Watson process with promiscuous mating: extinction probabilities in the supercritical case”. *Ann. Appl. Probab.* 6:922–939 (1996).
2. Alsmeyer, G., and Rösler, U. “Asexual versus promiscuous bisexual Galton-Watson processes: The extinction probability ratio”. *Ann. Appl. Probab.* 12:125–142 (2002).
3. Bruss, F.T. “A note on extinction criteria for bisexual Galton-Watson processes”. *J. Appl. Probab.* 21:915–919 (1984).
4. Daley, D. J. “Extinction conditions for certain bisexual Galton-Watson branching processes”. *Z. Wahrscheinlichkeith.* 9:315–322 (1968).
5. Daley, D.J., Hull, D.M., and Taylor, J.M. “Bisexual Galton-Watson branching processes with superadditive mating functions”. *J. Appl. Probab.* 23:585–600 (1986).
6. González, M., Molina, and M., Mota, M. “Limit behaviour for a subcritical bisexual Galton-Watson branching process with immigration”. *Statist. Probab. Lett.* 49:19–24 (2000).
7. González, M., Molina, M., and Mota, M. “On the limit behaviour of a supercritical bisexual Galton-Watson branching process with immigration of mating units”. *Stochastic Anal. Appl.* 19:933–943 (2001).
8. Ma, S. “Bisexual Galton-Watson processes in random environments”. *Acta Math. Appl. Sinica* 22:419-428 (2006).
9. Ma, S., Molina, M. “Two-sex branching processes with offspring and mating in a random environment”. *J. Appl. Probab.* 46:993-1004 (2009).
10. Ma, S., Xing, Y. “The asymptotic properties of supercritical bisexual Galton-Watson branching process with immigration of mating units”. *Acta Math. Sci.* 26:603-609 (2006).
11. Molina, M., Jacob, C., and Ramos, A. “Bisexual branching processes with offspring and mating depending on the number of couples in the population”. *Test* 17:245–281 (2008).
12. Molina, M., Mota M., and Ramos, A. “Bisexual Galton–Watson branching process with population–size dependent mating. *J. Appl. Probab.* 39:479–490 (2002).
13. Molina, M., Mota, M., and Ramos, A. “Bisexual Galton–Watson branching process in varying environments”. *Stochastic Anal. Appl.* 21:1353–1367 (2003).
14. Molina, M., Mota, M., and Ramos, A. “Limit behaviour for a supercritical bisexual Galton–Watson branching process with population–size dependent mating”. *Stochastic Proc. Appl.* 112:309–317 (2004).
15. Molina, M., Mota, M., and Ramos, A. “On L^α -convergence, $1 \leq \alpha \leq 2$, for a bisexual branching process with population–size dependend mating”. *Bernoulli* 12:457–468 (2006).
16. Molina, M., del Puerto, I., and Ramos, A. “A class of controlled bisexual branching processes with mating depending on the number of progenitor couples”. *Statist. Probab. Lett.* 77:1737–1743 (2007).
17. Müller, A., and Stoyan, D. *Comparison methods for stochastic models and risk.* John Wiley and sons, London (2002).
18. Xing, Y., and Wang, Y. “On the extinction of one class of population-size-dependent bisexual branching processes”. *J. Appl. Probab.* 42:175-184 (2005).

Nonparametric Inference in the Class of Controlled Two-sex Branching Models

Manuel Molina, Manuel Mota, and Alfonso Ramos

Department of Mathematics
University of Extremadura, Spain
(e-mails: mmolina@unex.es, mota@unex.es, aramos@unex.es)

Abstract. The class of two-sex branching models with random control on the number of progenitor couples is considered. For such a class, by considering that no assumptions are made about the functional form of the underlying offspring probability distribution, we obtain Bayes estimators for the offspring probability law and for its main moments. Also, we determine the corresponding 95% highest posterior density credibility sets. By way of illustration, we present some simulated examples where we check the accuracy of both the estimates and their corresponding 95% highest posterior density credibility sets.

Keywords: Branching models, Two-sex models, Controlled models, Nonparametric inference, Bayesian inference.

1 Introduction

Inside the general context of stochastic modelling, the branching process theory provides mathematical models to describe the probabilistic evolution of systems whose components (cells, particles, individuals in general) after certain life period reproduce and die. It is an active research area of both theoretical interest and applicability to such fields as biology, demography, ecology, epidemiology, genetics, medicine, population dynamics, and physics. Some classical monographs about this theory are Asmussen and Hering[2], Athreya and Ney[3], Guttorp[8] and Harris[10]. From an applied point of view one may cite the books by Jagers[12], Kimmel and Axelrod[13], Pakes[21] and Haccou *et al.*[9] which include practical applications to cell kinetics, cell biology, chemotherapy, gene amplification, human evolution, and molecular biology.

In particular, with the purpose to model the probabilistic evolution of populations where females and males coexist and form couples (femalemale), several classes of discrete time two-sex branching models have been studied. They include the bisexual Galton-Watson model (see Alsmeyer and Rösler[1], Bruss[4], Daley[5], Daley *et al.*[6]), models with immigration (see Gonzalez *et al.*[7], Ma and Xing[15]), in varying or in random environments (see Molina *et al.*[18], Ma and Molina[14]) and those models depending on the number of couples in the population (see Molina *et al.*[17], Xing and Wang[22]). We refer the reader to Hull[11] or Haccou *et al.*[9] for surveys of two-sex branching

models. However, the range of processes studied is not large enough in order to get an optimum modelling in many two-sex populations where a control on the number of couples in the population is required. It can be stated that significant efforts have been made regarding random control branching models with asexual reproduction. Now similar efforts should be made to develop models with a random control where reproduction is bisexual. We consider a class of controlled two-sex models where, in each generation, a random control on the number of couples that take part in the reproduction (progenitor couples) introduced in Molina *et al.*[19].

The paper is organized as follows: In the Section 2, the controlled two-sex model is described formally and interpreted intuitively. In Section 3, considering that no assumptions are made about the functional form of the underlying offspring distribution, we provide some results about the Bayesian estimation concerning the offspring law and its main moments. We also determine the corresponding 95% highest posterior density credibility sets.

2 The controlled two-sex model

The controlled two-sex branching process $\{(F_n, M_n)\}_{n \geq 1}$ is defined in the following form:

$$(F_{n+1}, M_{n+1}) = \sum_{i=1}^{\phi_{Z_n}} (f_{n,i}, m_{n,i}), \quad Z_{n+1} = L_{Z_n}(F_{n+1}, M_{n+1}), \quad n \in Z^+ \quad (1)$$

where the empty sum is considered to be $(0, 0)$ and Z^+ denotes the set of nonnegative integers. The random vector (F_{n+1}, M_{n+1}) represents the number of females and males in the $(n + 1)$ th generation. These females and males form Z_{n+1} couples. Initially, we assume that there are a positive number k_0 of couples in the population, i.e. $Z_0 = k_0$. The random vectors $\{(f_{n,i}, m_{n,i})\}_{n \geq 0; i \geq 1}$ are nonnegative, independent and identically distributed. Intuitively, $(f_{n,i}, m_{n,i})$ represents the number of females and males descending from the i th couple of the n th generation. If, for some positive integer n and $k \in Z^+$, $Z_n = k$, then:

- (a) L_k is a mating function. It is defined on $R^+ \times R^+$ and taking values in R^+ , where R^+ is the set of nonnegative real numbers, and it is assumed to be nondecreasing in each argument, integer-valued on the integers, and such that, for $x, y \in R^+$, $L_k(x, 0) = L_k(0, y) = 0$.
- (b) ϕ_k is a nonnegative integer-valued random control variable. The role of ϕ_k is to control the number of couples which will take part in the reproduction process. In fact, if $\phi_k > k$ then $\phi_k - k$ new couples are introduced in the population; if $\phi_k < k$ then $k - \phi_k$ couples leave the population and consequently, they do not participate in the reproduction

process; and no control is made if $\phi_k = k$. We will assume that $P(\phi_0 = 0) = 1$.

It then follows that in addition to its theoretical interest, this class of two-sex models also has clear practical implications, especially in population dynamics. For certain sexually reproducing animal population, it is reasonable to assume that the number of progenitor couples could be affected, in each generation, by random factors as weather conditions, food supply, fertility parameters, and so on. For example, in making policy decisions as to whether to introduce or re-introduce certain animal species into an environment, this class of models may provide appropriate mathematical models with which to describe the probabilistic behaviour of the population. Indeed, the motivation behind the class of models presented in 1 is the interest in developing two-sex models for such phenomena. As particular case, it includes the two-sex models introduced by Daley[5] and by Molina *et al.*[17], and generalizes to random control setting the model considered by Molina *et al.*[20].

3 Nonparametric estimation

We now consider a controlled two-sex branching process such that no assumption is made about the functional form of the underlying offspring distribution, so we consider a nonparametric setting. By simplicity, such a distribution will be denoted as $p = (p_{k,l} = P(f_{0,1} = k, m_{0,1} = l) : (k, l) \in S)$ where the support $S = \{(k, l) \in Z^+ \times Z^+ : p_{k,l} > 0\}$ is a finite set.

Let be $(\mu_1, \mu_2) = E[(f_{0,1}, m_{0,1})]$ and $(\sigma_{ij})_{i,j=1,2} = Cov[(f_{0,1}, m_{0,1})]$, the offspring mean vector and covariance matrix, respectively. We assume that $\sigma_{ij} < \infty, i, j = 1, 2$.

We shall assume the observation of the entire family tree, up to the n -th generation, namely $\{\phi_{Z_i}, (f_{i,j}, m_{i,j}); i = 0, \dots, n, j = 1, \dots, \phi_{Z_i}\}$. Let us denote by

$$Z_{i,(k,l)} = \sum_{j=1}^{\phi_{Z_i}} \mathbf{1}_{\{(f_{i,j}, m_{i,j})=(k,l)\}}, \quad (k, l) \in S$$

the number of couples in the i -th generation giving rise to exactly k females and l males. It is clear that

$$Z_i = \sum_{(k,l) \in S} Z_{i,(k,l)} \quad \text{and} \quad (F_{i+1}, M_{i+1}) = \sum_{(k,l) \in S} (k, l) Z_{i,(k,l)}.$$

It is easy to verify that the likelihood function satisfies

$$\ell(p) \propto \prod_{(k,l) \in S} p_{k,l}^{Y_{n,(k,l)}} \tag{2}$$

where $Y_{n,(k,l)} = \sum_{i=0}^n Z_{i,(k,l)}$ represents the total number of couples in the first n generations which have produced exactly k females and l males.

Considering (2), an appropriate conjugate class of prior distributions is the Dirichlet family,

$$\pi(p) = D_\tau \prod_{(k,l) \in S} p_{k,l}^{\tau_{k,l}-1} \quad (3)$$

where $\tau = (\tau_{k,l} : (k,l) \in S)$, $\tau_{k,l} > 0$, $D_\tau = \prod_{(k,l) \in S} \Gamma(\tau_{k,l})^{-1} \Gamma(\tau_*)$ and $\tau_* = \sum_{(k,l) \in S} \tau_{k,l}$. We refer the reader to Mendoza and Gutierrez-Peña[16] where some comments about the convenience of this class of distributions are given and some methods for deriving noninformative priors, including Jeffrey's rule, reference analysis and vague priors, are discussed.

Denoting by $\mathcal{F}_n^* = \sigma((f_{i,j}, m_{i,j}), i = 0, \dots, n; j = 1, \dots, \phi_{Z_i})$ and taking into account (2) and (3), we deduce that the posterior distribution is the Dirichlet law,

$$\pi(p|\mathcal{F}_n^*) = D_\gamma \prod_{(k,l) \in S} p_{k,l}^{\gamma_{k,l}-1}$$

with vector of parameters $\gamma = (\gamma_{k,l} : (k,l) \in S)$ where $\gamma_{k,l} = \tau_{k,l} + Y_{n,(k,l)}$. In particular, the marginal posterior distribution of $p_{k,l}$ is a Beta law with parameters $\gamma_{k,l}$ and $\gamma_* - \gamma_{k,l}$, where $\gamma_* = \sum_{(k,l) \in S} \gamma_{k,l} = \tau_* + \sum_{i=0}^n Z_i$. Assuming squared error loss function, we obtain the Bayes estimator for $p_{k,l}$,

$$\tilde{p}_{k,l} = E[p_{k,l} | \mathcal{F}_n^*] = \gamma_*^{-1} \gamma_{k,l} = \left(\sum_{i=0}^n Z_i + \tau_* \right)^{-1} (\tau_{k,l} + Y_{n,(k,l)}). \quad (4)$$

Next result provides the Bayes estimators of the offspring mean vector and the offspring covariance matrix.

Proposition 1. *Given a controlled two-sex branching model, the Bayes estimators of μ_i and σ_{ij} , $i, j = 1, 2$, under squared error loss function and assuming the class of conjugate prior distributions given in (3), are:*

(i)

$$\tilde{\mu}_i = \gamma_*^{-1} \sum_{(k_1, k_2) \in S} k_i \gamma_{k_1, k_2}, \quad i = 1, 2.$$

(ii)

$$\tilde{\sigma}_{ij} = \frac{1}{\gamma_*(1 + \gamma_*)} \left(\gamma_* \sum_{(k_1, k_2) \in S} k_i k_j \gamma_{k_1, k_2} - \sum_{(k_1, k_2), (l_1, l_2) \in S} k_i l_j \gamma_{k_1, k_2} \gamma_{l_1, l_2} \right)$$

$i, j = 1, 2$.

Proof.

(i) Using (4), we have for $i = 1, 2$,

$$\begin{aligned} \tilde{\mu}_i &= E \left[\sum_{(k_1, k_2) \in S} k_i p_{k_1, k_2} \mid \mathcal{F}_n^* \right] = \sum_{(k_1, k_2) \in S} k_i E[p_{k_1, k_2} \mid \mathcal{F}_n^*] \\ &= \gamma_*^{-1} \sum_{(k_1, k_2) \in S} k_i \gamma_{k_1, k_2}. \end{aligned}$$

(ii) For $i, j = 1, 2$,

$$\begin{aligned} \tilde{\sigma}_{ij} &= E \left[\sum_{(k_1, k_2) \in S} (k_i - \mu_i)(k_j - \mu_j) p_{k_1, k_2} \mid \mathcal{F}_n^* \right] \\ &= \sum_{(k_1, k_2) \in S} k_i k_j E[p_{k_1, k_2} \mid \mathcal{F}_n^*] - \sum_{(k_1, k_2) \in S} k_i k_j E[p_{k_1, k_2}^2 \mid \mathcal{F}_n^*] \\ &\quad - \sum_{(k_1, k_2) \neq (l_1, l_2)} k_i l_j E[p_{k_1, k_2} p_{l_1, l_2} \mid \mathcal{F}_n^*]. \end{aligned}$$

Using the fact that

$$E[p_{k_1, k_2}^2 \mid \mathcal{F}_n^*] = (\gamma_* (\gamma_* + 1))^{-1} \gamma_{k_1, k_2} (\gamma_{k_1, k_2} + 1)$$

and

$$E[p_{k_1, k_2} p_{l_1, l_2} \mid \mathcal{F}_n^*] = (\gamma_* (\gamma_* + 1))^{-1} \gamma_{k_1, k_2} \gamma_{l_1, l_2}, \quad (k_1, k_2) \neq (l_1, l_2),$$

the proof is completed.

Using the posterior distribution $\pi(\theta_1, \theta_2 \mid \mathcal{F}_n)$ we can determine sets of probable values of (θ_1, θ_2) . The most common procedure is based on looking at the points where the posterior density takes the highest values, namely $I(c) = \{(\theta_1, \theta_2) : \pi(\theta_1, \theta_2 \mid \mathcal{F}_n) \geq c\}$ where the constant c is chosen such that, given a credibility coefficient $1 - \alpha$,

$$\int_{I(c)} \pi(\theta_1, \theta_2 \mid \mathcal{F}_n) d\theta_1 d\theta_2 = 1 - \alpha.$$

We say that $I(c)$ is a high posterior density (HPD) credibility set.

In particular, from the posterior marginal densities of μ_i and σ_{ij} , $i, j = 1, 2$ we could derive HPD credibility sets. However, it is not easy to compute the posterior densities for such parameters. In these cases, Monte Carlo approximations can be calculated, by simulating a sufficiently large number of values for (θ_1, θ_2) according to the posterior density $\pi(\theta_1, \theta_2 \mid \mathcal{F}_n)$.

References

1. Alsmeyer, G, and Rösler, U., “The bisexual Galton-Watson process with promiscuous mating: extinction probabilities in the supercritical case”, *Annals of Applied Probability* 6, 922-939 (1996).
2. Asmussen, G., and Hering, H., *Branching Processes*, Birkhäuser (1983).
3. Athreya, K., and Ney, P., *Branching Processes*, Springer-Verlag (1972).
4. Bruss, F. T., “A note on extinction criteria for bisexual GaltonWatson processes”, *Journal of Applied Probability* 21, 915-919 (1984).
5. Daley, D. J., “Extinction conditions for certain bisexual Galton-Watson branching processes”, *Zeitschrift für Wahrscheinlichkeitstheorie* 9, 315-322 (1968).
6. Daley, D. J., Hull, D. M., and Taylor, J. M., “Bisexual Galton-Watson branching processes with superadditive mating functions”, *Journal of Applied Probability* 23, 585-600 (1986).
7. González, M., Molina, M., and Mota, M., “Limit behaviour for a subcritical bisexual Galton- Watson branching process with immigration”, *Statistics and Probability Letters* 49, 19-24 (2000).
8. Guttorp, P., *Statistical Inference for Branching Processes*, Wiley (1991).
9. Haccou, P., Jagers, P., and Vatutin, V., *Branching Processes: Variation, Growth, and Extinction of Populations*, Cambridge University Press (2005).
10. Harris, T., *The Theory of Branching Processes*, Springer-Verlag (1963).
11. Hull, D. M., “A survey of the literature associated with the bisexual Galton-Watson branching process”, *Extracta Mathematicae* 18, 321-343 (2003).
12. Jagers, P., *Branching Processes with Biological Applications*, Wiley (1975).
13. Kimmel, M., and Axelrod, D. E., *Branching Processes in Biology*, Springer-Verlag (2002).
14. Ma, S., and Molina, M., “Two-sex branching processes with offspring and mating in a random environment”, *Journal of Applied Probability* 46, 993-1004 (2009).
15. Ma, S., and Xing, Y., “The asymptotic properties of supercritical bisexual GaltonWatson branching processes with immigration of mating units”, *Acta Mathematicae Sciences* 26, 603-609 (2006).
16. Mendoza, M., and Gutiérrez-Peña, E., “Bayesian conjugate analysis of the Galton–Watson process”, *Test* 9, 149–172 (2000).
17. Molina, M., Mota, M., and Ramos, A., “Bisexual GaltonWatson branching process with populationsize dependent mating”, *Journal of Applied Probability* 39, 479-490 (2002).
18. Molina, M., Mota, M., and Ramos, A., “Bisexual GaltonWatson branching process in varying environments”, *Stochastic Analysis and Applications* 21, 1353-1367 (2003).
19. Molina, M., Mota, M., and Ramos, A., “Two-sex branching models with random control on the number of progenitor couples”, *To appear in Methodology and Computing in Applied probability*.
20. Molina, M., del Puerto, I., and Ramos, A., “A class of controlled bisexual branching processes with mating depending on the number of progenitor couples”, *Statistics and Probability Letters* 77, 1737-1743 (2007).
21. Pakes, A., *Biological Applications of Branching Processes*, Handbook of Statistics, v. 21, C.N. Shanbhag and C.R. Rao, Eds, Elsevier Sciences B.V. (2003).

22. Xing, Y., and Wang, Y., “On the extinction of one class of population-size-dependent bisexual branching processes”, *Journal of Applied Probability* 42, 175–184 (2005).

Factor Analysis (FA) as ranking and an Efficient Data Reducing approach for decision making units

Reza Nadimi*, Hamed Shakouri G. *, Reza Omid†

*Department of Industrial Engineering, Faculty of Engineering, University of Tehran, Tehran, Iran.

†Department of Social Policy, Faculty of Social Science, University of Tehran, Tehran, Iran.

Corresponding author: h.shakouri@gmail.com

Abstract- This article compares two techniques: Data Envelopment Analysis (DEA) and Factor Analysis (FA) to aggregate multiple inputs and outputs in the evaluation of decision making units (DMU). Data envelopment analysis (DEA), a popular linear programming technique, is useful to rate comparatively operational efficiency of DMUs based on their deterministic or stochastic input-output data. Factor analysis techniques, such as Principal Components Analysis, have been proposed as data reduction and classification technique, which can be applied to evaluate of decision making units (DMUs). FA, as a multivariate statistical method, combines new multiple measures defined by inputs/outputs. Nonparametric statistical tests are employed to validate the consistency between the ranking obtained from DEA and FA. Also, the results have been compared with PCA approach. Results of numerical reveal that new approach shows a consistency in ranking with DEA.

Keywords: Decision Making; Data Envelopment Analysis; Factor Analysis, Principal Component Analysis.

I- Introduction

This article proposes a Factor Analysis (FA) approach to evaluate of decision making units (DMUs). In this method, FA is used as a new approach to ranking of decision making units and data reduction. Moreover, correlation between rankings obtained by FA and DEA techniques is much higher than what is gained from the PCA&DEA method, which is introduced by Zhu [2].

The rest of this article is organized as follows. In Section 2, a brief description of the DEA models used for ranking of DMUs is presented. Section 3 gives the fundamental of FA technique. The FA approach is developed in Section 4. Numerical comparison of the proposed FA method versus DEA and PCA procedures is presented in Section 5, using several benchmark data to evaluate consistency of each method. Finally, Section 6 concludes this research.

II- Data Envelopment Analysis

Data envelopment analysis (DEA), is analytical tool which first introduced by Charnes et al.[1], in 1978. It is the performance measurement technique that applies to evaluation the relative efficiency of decision-making units (DMU's) in organization such as banks, dental services, police, motor registries, hospitals etc.

Various models, used for ranking of DMUs, such as CCR [1], BCC [3] and ADD [4] are applied. The standard DEA method assigns an efficiency score less than one to inefficient DMUs, from which a ranking can be derived. However, efficient DMUs all have an efficiency of 1, so that for these units no ranking can be given. Andersen and Petersen (AP model) achieve a full ranking by undertaking a DEA without assessing the DMU itself[5]. In fact, they proposed the idea of modifying the envelopment LP formulation so that the corresponding column of the DMU being scored is removed from the coefficients matrix. Thus we use the AP-model as a basis to rank the relative efficiency of DMUs with unit efficiency, in order to compare validity of other assessment techniques in this paper. AP-model, (1), can be written as follows:

$$J_k^* = \min w_k$$

subject to

$$\sum_{\substack{i=1 \\ i \neq k}}^m \alpha_i \mathbf{x}_i \leq \mathbf{x}_k w_k, \quad (1)$$

$$\sum_{\substack{i=1 \\ i \neq k}}^m \alpha_i \mathbf{y}_i \geq \mathbf{y}_k,$$

$$\alpha_i \geq 0, \quad i = 1, \dots, m$$

The program depends on evaluating the k^{th} unit; where $\mathbf{x}_k = [x_{1j}, x_{2j}, \dots, x_{mj}]$, and $\mathbf{y}_k = [y_{1j}, y_{2j}, \dots, y_{sj}]$, denote the nonnegative vector of input and output values for DMU_k respectively. Hence, each J_k^* lies between 0 and $+\infty$. Also, In model (1), α_j is the *Factor weights*.

However, the super-efficient methodology can give "specialized" DMUs an excessively high ranking.

Consequently, in this paper we apply the *Factor Analysis* (FA) to reduce data; indeed, we use this method to evaluate and rank DMUs while minimizing loss of the information.

III- Factor Analysis (FA)

Factor Analysis is a statistical method that is based on the correlation analysis of multi-variables. The main applications of factor analytic techniques are: (1) to reduce the number of variables and (2) to detect structure in the relationships between variables, in order to classify variables. Therefore, factor analysis is applied as a data reduction or structure detection method.

It can be used as a method to data reduction. R. Nadimi, F. Jolai[12] applied combination of factor analysis and data envelopment analysis to data reduction in decision making units. They used factor analysis as a method to lessen the number of data. In follow, data envelopment analysis was used with combination of factor analysis to data ranking. But in this paper factor analysis only is used as a new method in ranking of data.

There are two major types of FA: *exploratory* and *confirmatory*. In exploratory FA, one seeks to describe and summarize data by grouping together variables that are correlated. The variables themselves may or may not have been chosen with potential underlying processes in mind. Exploratory FA is usually performed in the early stages of research, when it provides a tool for consolidating variables and for generating hypotheses about underlying processes. Confirmatory FA is a much more sophisticated technique used in the advanced stages of the research process to test a theory about latent processes. Variables are carefully and specifically chosen to reveal underlying processes [6].

To explain the method, a few terms are defined. The first terms involve correlation matrices. The correlation matrix produced by the observed variables is called the *observed correlation matrix*. The correlation matrix produced from factors is called the *reproduced correlation matrix*. The difference between observed and reproduced correlation matrices is called *residual correlation matrix*. In a good FA, correlations in the residual matrix are small, indicating a close fit between the observed and reproduced matrices [6]. Then, factors are formed by grouping the variables that have higher correlation with each other.

Let $\mathbf{d}_{(n \times 1)}$ be a random vector with a mean of $\boldsymbol{\mu}$ and a covariance matrix named $\boldsymbol{\Sigma}_{(p \times p)}$, where d_i specifies efficiency or an overall performance index of the i^{th} DMU. Then a k -factor model holds for \mathbf{d} , if it can be written in the following form:

$$\mathbf{d} = \mathbf{H}\mathbf{f} + \mathbf{u} + \boldsymbol{\mu} \quad (2),$$

where $\mathbf{H}_{(n \times k)}$ is a matrix of constants and $\mathbf{f}_{(k \times 1)}$ and $\mathbf{u}_{(n \times 1)}$ are random vectors. The elements of \mathbf{f} are called common

factors and the elements of \mathbf{u} are *specific* or *unique* factors. In this study we shall suppose that:

$$\begin{aligned} E(\mathbf{f}) &= \mathbf{0}, \text{Cov}(\mathbf{f}) = \mathbf{I} \\ E(\mathbf{u}) &= \mathbf{0}, \text{Cov}(u_i, u_j) = 0; i \neq j \\ \text{Cov}(\mathbf{f}, \mathbf{u}) &= \mathbf{0} \end{aligned} \quad (3).$$

Thus, if (2) holds, the covariance matrix of \mathbf{d} can be split into two parts, as follows:

$$\boldsymbol{\Sigma} = \mathbf{H}\mathbf{H}^T + \boldsymbol{\Phi} \quad (4),$$

where $\mathbf{H}\mathbf{H}^T$ is called the *communality* and represents the variance of q_i which is shared with the other variables via the common factors and $\boldsymbol{\Phi} = \text{Cov}(\mathbf{u})$ is called the *specific* or *unique variance* and is due to the unique factors \mathbf{u} . This matrix explains the variability in each q_i that is not shared with the other variables. The main goal of FA is to apply \mathbf{f} instead of \mathbf{d} for assessing DMUs. To do this, mainly there are three main stages in a typical FA technique [7]:

1. Initial solution: Variables, as indexes of DMU performance measures, are selected and an inter-correlation matrix is generated. An inter-correlation matrix is a $p \times p$ array of the correlation coefficients of p variables with each other. Usually, each variable is standardized by a certain formula, e.g. to have a mean of 0.0 and a standard deviation of 1.0. When the degree of correlation between the variables is weak, it is not feasible for these variables to have a common factor, and a correlation between these variables is not studied. Kaiser–Meyer–Olkin (KMO) and Bartlett’s tests of sphericity (BTS) are then applied to the studied variables in order to validate if the remaining variables are factorable.

2. Extracting the factors: An appropriate number of components (Factors) are extracted from the inter-correlation matrix based on the initial solution. Due to the standardization method, there should be a certain rule to extract the selected effective factors.

3. Rotating the factors: Sometimes one or more variables may load about the same on more than one factor, making the interpretation of the factors ambiguous. Thus, factors are rotated in order to clarify the relationship between the variables and the factors. While various methods can be used for factor rotation, the Varimax method is the most commonly used one.

Let’s summarize and formulize the above steps as follows. In this study, we skip the rotation step.

First, the correlation matrix, namely \mathbf{R} , is computed on the basis of data due to the standardized variables, d_{ij} :

$$\mathbf{R} = \text{Corr}(\mathbf{D}) = \mathbf{D}^T \mathbf{D} \quad (5),$$

where, \mathbf{D} is an $n \times p$ matrix of p variables for n DMU’s.

This matrix can be decomposed to a product of three matrices:

$$\mathbf{R} = \mathbf{V}\mathbf{L}\mathbf{V}^T \quad (6),$$

where, \mathbf{V} is the $p \times p$ matrix of eigenvectors and $\mathbf{L} = \text{Diag}([\lambda_1, \dots, \lambda_p])$ is a diagonal matrix of the eigenvalues, assorted descendingly.

At the second step, different criteria may be applied to extract the most important factors. Since sum of the first r eigenvalues divided by the sum of all the eigenvalues, $(\lambda_1+\lambda_2+\dots+\lambda_r)/(\lambda_1+\lambda_2+\dots+\lambda_p)$, represents the “proportion of total variation” explained by the first r factor components, we select r principal components as the factors, if $(\lambda_1+\lambda_2+\dots+\lambda_r)/(\lambda_1+\lambda_2+\dots+\lambda_p) > 90\%$. Another criterion is to cut the matrix \mathbf{L} from a point that the ratio of λ_i/λ_{i+1} is maximized. However, r eigenvalues are defined as dominant eigenvalues. The dominant eigenvalues are saved and the other are skipped. To explain more, suppose \mathbf{L} and \mathbf{V} are decomposed as follows:

$$\mathbf{L} = \begin{bmatrix} \mathbf{L}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{L}_2 \end{bmatrix} \quad (7)$$

where $\mathbf{L}_1 (r \times r)$ and \mathbf{L}_2 are diagonal matrixes. Consequently, the eigenvectors \mathbf{V} will be separated into two parts too:

$$\mathbf{V} = [\mathbf{V}_1, \mathbf{V}_2] \quad (8)$$

Similarly, \mathbf{V}_1 and \mathbf{V}_2 are $p \times r$ and $p \times (p-r)$ matrices, respectively. Suppose (6) is rewritten as follows:

$$\mathbf{R} = (\mathbf{V} \sqrt{\mathbf{L}}) (\sqrt{\mathbf{L}} \mathbf{V}^T) \quad (9)$$

Then, replacing \mathbf{L} with the form given by (7), the first part $\mathbf{V}_1 \sqrt{\mathbf{L}_1}$ is called the *Factor Loading* matrix and denoted by $\mathbf{A} (p \times r)$. Equation (9) is frequently called the fundamental equation for FA. It represents the assertion that the correlation matrix is a product of the factor loading matrix, \mathbf{A} , and its transpose [6]. It can be shown that an estimate of the unique or specific variance matrix, Φ , in (4) is:

$$\mathbf{B} = \mathbf{I} - \mathbf{A} \mathbf{A}^T \quad (10)$$

where $\mathbf{I} (p \times p)$ is the identity matrix.

So far our study of the factor model has been concerned with the way in which the observed variables are functions of the (unknown) factors, \mathbf{f} . Instead, factor scores can be estimated by the following pseudo-inverse method:

$$\mathbf{S}^T = (\mathbf{A}^T \mathbf{B}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{B}^{-1} \quad (11)$$

$$\mathbf{F} = \mathbf{D} \mathbf{S} \quad (12)$$

where \mathbf{F} is a $n \times r$ matrix, each row of which corresponds to a DMU. The estimate in (12) is known as Bartlett’s factor score, and \mathbf{S} is called the *factor score* coefficient matrix.

In this paper, we use the FA technique to evaluate DMUs by reducing inputs and outputs whilst minimizing the loss of information. This will be introduced in the next section.

IV- New approach: FA method

In here, ratios of individual output to individual input is used to describe of proposed approach. Thus this proportion is applied to evaluate and rank DMUs according to their performances which are given as follows:

$$d_{ir}^j = y_{rj} / x_{ij} ; i=1, \dots, m; r=1, \dots, s; \quad (13)$$

$j=1, \dots, n$

for each DMU $_j$. Where the d_{ir}^j gives the ratio between every output and every input. Obviously, the bigger the d_{ir}^j , the better the performance of DMU $_j$ in terms of the r^{th} output and the i^{th} input [8].

Now let $d_k^j = d_{ir}^j$, with, e.g. $k=1$ corresponds to $i=1, r=1$ and $k=2$ corresponds to $i=1, r=2$, etc., where $k=1, \dots, p'$; $p'=m \times s$; for example: $d_1 = y_1/x_1, d_2 = y_1/x_2, \dots$

We need to find some weights that combine those p' individual ratios of d_k^j for DMU $_j$. Consider the following

$n \times p'$ data matrix, composed by d_k^j ’s: $\mathbf{D}^T = [d_1, \dots, d_{p'}]_{n \times p'}$,

where each row represents p' individual ratios of d_k^j for each DMU and each column represents a specific output/input ratio, i.e. $\mathbf{d}_k = [d_k^1, \dots, d_k^n]^T$. In a modified approach, proposed by Premachandra [9], \mathbf{D}^T is re-defined as an augmented matrix, the ending column of which is equivalent to the sum of the elements in the first p' columns of the original matrix:

$$d_{p'+1}^j = \sum_{k=1}^{p'} d_k^j \quad j=1, \dots, n \quad (14)$$

The new added variable, is supposed to take into account the overall performance of each DMU with respect to all the variables d_{ir}^j . As a normalizing skill, each column is then divided by its least element, thus a new matrix, $\mathbf{D} (p \times n)$; $p=p'+1$, is generated which will be processed from now on.

In this paper, the factor analysis is employed to find out new independent measures which are respectively different linear combinations of d_1, \dots, d_p . In fact, we apply the estimation given by (12) to obtain factor scores, thus, the FA process of \mathbf{D} is carried out as follows:

Step 1: Calculate the sample correlation matrix, given by (5), to obtain eigenvalues and eigenvectors (solutions to $|\mathbf{R} - \lambda \mathbf{1}_p| = 0$ where $\mathbf{1}_p$ is a $p \times p$ identity matrix), as introduced in (6).

Step 2: Considering $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ as the sorted eigenvalues, compute the following weightings, which determine share of each factor in the model:

$$w_i = \frac{\lambda_i}{\sum_{k=1}^p \lambda_k} ; i=1, \dots, p \quad (15)$$

Each weighting actually determines the share of each eigenvalue out of a whole. This approach uses the same method of Zhu [3] to obtain sign of the weightings w_i , i.e. if

sum of the corresponding eigenvector elements is positive, then w_i is considered positive, otherwise it is negative.

Step 3: Apply FA technique on D to obtain S^F and then F , as defined by (11) and (14).

Step 4: Select the factor components by determination of the dominant eigenvalues according to one of the criteria proposed in Section 3.

Step 5: Compute:

$$z = \sum_{i=1}^r w_i f_i \quad (16),$$

where f_i is the i^{th} column of the matrix F in (14) and r is the number of the dominant eigenvalues. The value of z gives a combined measure to evaluate and rank performance of DMUs.

V- Numerical results

The proposed method is applied to several sets of sample data, the numerical results of which are illustrated and compared to other methods in this section.

Example1: In this example, we apply data set used by Wong et al. [11], to compare efficiencies of seven university departments. Three inputs and three outputs are defined as follows, data of which is listed in Table 1.

x_1 : Number of academic staff

x_2 : Academic staff salaries

x_3 : Support of undergraduate students

y_1 : Number of undergraduate students

y_2 : Number of postgraduate students

y_3 : Number of research papers published

[11]: Data set used by Wong et al.1Table

DMU	x_1	x_2	x_3	y_1	y_2	y_3
dmu1	12	400	20	60	35	17
dmu2	19	750	70	139	41	40
dmu3	42	1500	70	225	68	75
dmu4	15	600	100	90	12	17
dmu5	45	2000	250	253	145	130
dmu6	19	730	50	132	45	45
dmu7	41	2350	600	305	159	97

The same procedure of section 4 is followed. The matrix D is generated by 10 variables extracted out of data in Table 1, and four dominant eigenvectors are selected. Table 2 illustrates eigen-analysis applied for PCA and FA, and Table 3 includes the results of ranking.

: Eigen-analysis for FA and PCA 2Table approaches

Eigen values	4.15	3.09	1.73	0.85
Shares of Eigen values (w_i)	0.41	0.31	0.17	-0.08
Eigen vector	v_1	v_2	v_3	v_4
v_{i1}	-0.08	0.33	0.53	-0.40
v_{i2}	0.11	-0.23	0.63	-0.26
v_{i3}	0.33	-0.41	0.09	-0.02
v_{i4}	0.24	0.43	-0.25	-0.28
v_{i5}	0.39	0.20	-0.29	-0.29
v_{i6}	0.38	-0.28	-0.24	-0.23

In this example the correlation between results obtained by PCA (Zhu) and DEA is 0.321, while correlation between DEA&PCA (PM) is 0.678. However, the new approach of FA riches to a higher correlation with the DEA, that is 0.75, due to the scores given to the dmu5 and dmu6. This example shows that the FA approach can lead to better results, in the sense of DEA ranking, compared to the both PCA approaches proposed by Zhu and Premachandra.

Example 2: As the last case, we compared the PCA (Zhu), PCA (PM) and FA approaches on the base of the DEA approach as performed in Kim et al. [10] for 33 telephone offices in S. Korea (See Table 4 for more information). Corresponding correlations which are given in Table 5, are 0.63, 0.75, and 0.77 respectively. While all the methods are statistically significant at 1% level, the new method based on FA shows better capability for ranking.

VI- CONCLUSION

The current article presents alternative approach to rank and evaluate DMUs which have multiple outputs and multiple inputs. The DEA –non-statistical method– uses linear programming technique to obtain a ratio between weighted outputs and weighted inputs. The new approach proposed in this paper is applied to evaluate efficiencies and rank DMUs. Factor analysis is a multivariate statistical method that uses information obtained from eigenvalues to reduce data. Results obtained by numerical experiments employed, show that there is a high correlation between DEA and FA methods, even higher than what obtained by the PCA methods. Thus, we can use FA to evaluate efficiency and ranking DMUs instead of DEA with significance and minimum lose of information.

: **Efficiencies and rankings obtained by the three methods**3Table

DMU	DEA		PCA(Zhu)		PCA(PM)		FA(New method)	
	Score	Rank	Score	Rank	Score	Rank	Score	Rank
dmu1	1.829615	1	0.51261	2	4.13838	1	2.011187	1
dmu2	1.048895	6	0.288772	4	3.315666	5	1.712316	5
dmu3	1.198308	4	0.011661	5	3.25405	6	1.559566	6
dmu4	0.819737	7	-1.9633	7	1.616393	7	0.895427	7
dmu5	1.219992	3	0.456634	3	3.801057	3	1.943119	2
dmu6	1.190642	5	0.918423	1	3.846452	2	1.917534	3
dmu7	1.266094	2	-0.2248	6	3.47953	4	1.883721	4

VII-REFERENCES

[1] Charnes, W.W. Cooper, E. Rhodes, Measuring the efficiency of decision making units, *European Journal of Operations Research* 2 (1978) 429 -444.

[2] J. Zhu, Data envelopment analysis vs principal component analysis: An illustrative study of economic performance of Chinese cities, *European Journal of Operation Research* 111,(1998) 50-61.

[3] R.D. Banker ,A. Charnes,W.W. Cooper,Some models for estimating technical and scale inefficiencies in data envelopment analysis, *Management Science* 30(9)(1984)1079-1092

[4] A. Charnes, W.W. Cooper,B. Golany,L. Seiford, Foundations of data envelopment analysis for Pareto-Koop-mans efficient empirical production functions, *Journals of Econometrics* 30 (1985) 91-107.

[5] P. Andersen, N.C. Petersen, A procedure for ranking efficient units in data envelopment analysis, *Management science* 39(10), (1993) 1261-1294.

[6] B.G. Tabachnick L.S. Fidell, using multivariate statistics, fourth edition, allyn& bacon person education company (1996) 582-627.

[7] How to perform and interpret Factor analysis using SPSS, www.ncl.ac.Uk/iss/statistics/docs/Factoranalysis.html, 2002.

[8] N. Adler, B. Golany, Evaluation of deregulated airline networks using data envelopment analysis combined with principal component analysis with an application to Western Europe, *European Journal of Operations Research* 132, (2001) 260-273.

[9] I.M. Premachandra, A note on DEA vs principal component analysis: An improvement to Joe Zhu’s approach, *Journal of Operational Research Society* 132(2001) 553-560.

[10]S.H. Kim, C.G. Park, K.S. Park, An application of Data Envelopment Analysis in telephone offices evaluation with partial data. *Computers & Operation Research* 26(1999), 59-72.

[11]Y.H.B. Wong, J.E. Beasley, Restricting weight flexibility in data envelopment analysis, *Journal of Operational Research Society* 41(9), (1990) 829-835.

[12]R. Nadimi, F. Jolai, “Joint Use of Factor Analysis (FA) and Data Envelopment Analysis (DEA) for Ranking of Data Envelopment Analysis”, *International Journal of Mathematical, Physical and Engineering Sciences* 2;4 © www.waset.org Fall 2008.

Table 4: Data for Telephone Office

DMU	X ₁	X ₂	X ₃	Y ₁	Y ₂	Y ₃	Y ₄	Y ₅
dmu1	239	7.03	158	47.1	16.67	34	28	2
dmu2	261	3.94	163	37.5	14.11	20	26	3
dmu3	170	2.1	90	20.7	6.8	12.6	19	3
dmu4	290	4.51	201	41.8	11.07	6.27	23	4
dmu5	200	3.99	140	33.4	9.81	6.49	30	2
dmu6	283	4.65	214	42.4	11.34	5.16	21	4
dmu7	286	6.54	197	47	14.62	13	9	2
dmu8	375	6.22	314	55.5	16.39	7.31	14	1
dmu9	301	4.82	257	49.2	16.15	6.33	8	3
dmu10	333	6.87	235	47.1	13.86	6.51	6	2
dmu11	346	6.46	244	49.4	15.88	8.87	18	2
dmu12	175	2.06	112	20.4	4.95	1.67	32	5
dmu13	217	4.11	131	29.4	11.39	4.38	33	2
dmu14	441	7.71	214	61.2	25.59	33	16	3
dmu15	204	3.64	163	32.3	9.57	3.65	15	4
dmu16	216	2.24	154	32.8	11.46	9.02	25	2
dmu17	347	5.65	301	59	17.82	8.19	29	1
dmu18	288	4.66	212	42.3	14.52	7.33	24	4
dmu19	185	3.37	178	33	9.46	2.91	7	2
dmu20	242	5.12	270	65.1	24.57	20.7	17	1
dmu21	234	2.52	126	31.6	8.55	7.27	27	2
dmu22	204	4.24	174	32.5	11.15	2.95	22	3
dmu23	356	7.95	299	66	22.25	14.9	13	2
dmu24	292	4.52	236	50	14.77	6.35	12	3
dmu25	141	5.21	63	21.5	9.76	16.3	11	2
dmu26	220	6.09	179	47.9	17.25	22.1	31	2
dmu27	298	3.44	225	42.4	11.14	4.25	4	2
dmu28	261	4.3	213	41.7	11.13	4.68	20	5
dmu29	216	3.86	156	31.6	11.89	10.5	3	3
dmu30	171	2.45	150	24.1	9.08	2.6	10	5
dmu31	123	1.72	61	12	4.78	2.95	5	1
dmu32	89	0.88	42	6.4	3.18	1.48	2	5
dmu33	109	1.35	57	10.6	3.43	2	1	4

: Efficiencies and rankings obtained by the three methods5Table

DMU	DEA		PCA(Zhu)		PCA(PM)		FA(New method)	
	Score	Rank	Score	Rank	Score	Rank	Score	Rank
dmu1	1.00000	3	2.11492	1	11.48890	1	0.67729	2
dmu2	1.00000	13	1.34052	7	9.39700	5	0.49045	6
dmu3	1.00000	11	1.48560	5	11.13670	2	0.75307	1
dmu4	0.86818	20	-0.51486	20	5.25150	17	-0.13271	16
dmu5	0.99367	18	0.49112	10	7.19350	11	0.19885	12
dmu6	0.84137	24	-0.69660	24	4.74490	22	-0.22910	21
dmu7	0.86995	29	-0.45802	18	4.52720	23	-0.27669	25
dmu8	0.72081	33	-1.04620	29	2.99440	31	-0.59070	32
dmu9	0.82025	26	-0.71586	25	3.44640	29	-0.38017	28
dmu10	0.75450	32	-1.33681	32	2.69360	33	-0.62922	33
dmu11	0.77697	31	-0.72447	26	4.01270	27	-0.40117	29
dmu12	1.00000	1	0.22024	12	9.27820	6	0.47063	8
dmu13	1.00000	12	0.46137	11	7.01790	13	0.15976	13
dmu14	1.00000	2	0.99676	8	7.62280	10	0.28591	10
dmu15	0.87213	23	-0.62439	21	4.96000	20	-0.16108	19
dmu16	1.00000	9	1.47139	6	8.31010	8	0.54230	4
dmu17	0.83311	7	-0.20785	14	4.50530	25	-0.25663	24
dmu18	0.84828	17	-0.18136	13	5.56040	14	-0.08199	14
dmu19	0.79771	30	-0.93114	28	3.23510	30	-0.44558	30
dmu20	1.00000	4	1.70960	4	7.08100	12	0.33851	9
dmu21	1.00000	16	0.89027	9	7.86410	9	0.47226	7
dmu22	0.84563	27	-0.43421	17	5.08270	18	-0.21236	20
dmu23	0.85252	14	-0.28974	15	4.25360	26	-0.29092	26
dmu24	0.89417	21	-0.47861	19	3.96960	28	-0.23178	23
dmu25	1.00000	10	1.78330	2	11.07270	3	0.67505	3
dmu26	1.00000	8	1.71207	3	9.54480	4	0.50040	5
dmu27	0.86546	22	-1.07657	30	2.70010	32	-0.48153	31
dmu28	0.88027	6	-0.63001	22	5.04150	19	-0.14681	17
dmu29	0.83361	28	-0.42143	16	4.92900	21	-0.23122	22
dmu30	0.91892	15	-0.63932	23	5.28470	16	-0.11320	15
dmu31	0.77394	25	-0.73120	27	4.51310	24	-0.35090	27
dmu32	1.00000	5	-1.11510	31	8.52320	7	0.23242	11
dmu33	0.93490	19	-1.42338	33	5.55320	15	-0.15314	18

A new concept in possibility of equality to create innovative constraints in fuzzy linear regression

R. Nadimi, S. F. Ghaderi

Department of Industrial Engineering, Faculty of Engineering, University of Tehran, Tehran, Iran.

Corresponding author: r.nadimi@yahoo.com

Abstract- Possibility of equality between two or more fuzzy numbers is a popular method to consider their degree of fitness. Possibility of equality may be applied to establish the constraints of fuzzy linear regression in which conjunction problem is under consideration. In this study, a new concept of the possibility of equality, that creates new restrictions, will be introduced and applied in fuzzy regression model, and then a more precise method will be represented to calculate the amount of error. To compare the performance of the proposed approach with the other methods, numerical examples are given. Total amount of error is calculated to confirm the efficiency of the proposed approach.

Keywords: Fuzzy linear regression; Possibility of equality; Fuzzy number.

I- INTRODUCTION

Regression analysis is a statistical method applied to consider the relationship between the dependent and independent variables. Fuzzy regression model is an extension of common regression in which one of the input and output data or both of them are regarded as fuzzy numbers. Probability distribution function is used to estimate parameters in classical regression and possibility theory which was introduced by [26] is applied to estimate fuzzy regression parameters.

Fuzzy linear regression was introduced by Tanaka et al. [25] in which the input and output data were crisp and fuzzy, respectively. It has been successfully implemented in several fields of forecasting ([21],[17],[18], [14],[2], [13], [10], [16], [3],[20],[19],[24]). In general, fuzzy regression models are classified into two classes:

- 1) The possibilistic model: Minimize the fuzziness of the model by minimizing the total spreads of its fuzzy coefficients, subject to covering the observed data by the estimated data of the model ([25], [23],[22]).
- 2) The least-squares model: Minimize the distance between estimated output of the model and the observed amount, based on their modes and spreads ([6], [11],[8],[4],[17]).

Fuzzy Number: Based on Dubois and Prade [7] \tilde{A} is defined a fuzzy number which satisfies the following criteria:

First: normality, $\exists x \in \mathbf{R}$ such that $\mu_{\tilde{A}}(x)=1$

Second: convexity, $\forall x_1, x_2 \in \mathbf{R}, \forall h \in [0;1]$

$\mu_{\tilde{A}}(hx_1+(1-h)x_2) \geq \min(\mu_{\tilde{A}}(x_1), \mu_{\tilde{A}}(x_2))$
 $\tilde{A} = (c_L, a, c_R)_{LR}$ is a LR-type fuzzy number where a, c_L and c_R are the center, left spread and right spread of fuzzy number, respectively ($c_L \& c_R > 0$). When $c_L=c_R=c$, we have a symmetric triangular fuzzy number. Thus, $\tilde{A} = (a, c)_L$ is a symmetric triangular fuzzy number if:

$$\mu_{\tilde{A}}(x) = L\left(\frac{a-x}{c}\right) = 1 - \frac{|a-x|}{c}, \quad a-c \leq x \leq a+c \quad (1)$$

In this paper symmetric triangular fuzzy numbers is only considered for simplicity.

Problem Definition: In order to define the possibility of equality between two fuzzy numbers, (Dubois and Prade [7]) proposed the following index:

$$Poss(\tilde{A}_1 = \tilde{A}_2) = \sup_{x \in \mathbf{R}^1} \min\{\mu_{\tilde{A}_1}(x), \mu_{\tilde{A}_2}(x)\} \quad (2)$$

where Poss is short for Possibility.

Finding out a suitable mathematical model along with the best fitting coefficients of the model from the observed data is one of the fuzzy regression analysis goals.

The *Min*, *Max* and *Conjunction* problems are three types of possibilistic linear regression analysis to gain the mentioned aim that dealt with by Tanaka et al. [24].

Conjunction problem is a popular method that uses the concept of fuzzy number inclusion to find the best fitting coefficients. Some papers evolve *conjunction* problem's constraints by the definition of the possibility equality of two fuzzy numbers (for more information see Shakouri G. and Nadimi [17],[12]). Besides that, Shakouri G. and Nadimi [17] introduced Non-equality possibility to designed a new objective function.

The scope of this paper is limited to the conjunction problem. It deals with some problems about the possibility equality index. Afterward, a new possibility equality index is presented to encounter the problems and to establish new constraints in fuzzy linear regression.

Figure 1 is considered based on the mentioned aim. \tilde{A}_1 is supposed a fuzzy number in a set of fuzzy number observations. Different objective functions will lead to different estimated parameters in which the above index is considered as a criterion to find out the best estimated coefficients. For instance, \tilde{A}_2 and \tilde{A}_3 can be two optional estimated fuzzy numbers for \tilde{A}_1 . The possibility equality

index is the same for both \tilde{A}_1 & \tilde{A}_2 and \tilde{A}_1 & \tilde{A}_3 but the common area is different for them. Meanwhile the spread of \tilde{A}_2 is wider than \tilde{A}_3 . Therefore estimating fuzzy linear regression parameters is caused a little deviation to actual amount of parameters considering the mentioned index. It drives us out of the fuzzy regression analysis aim. Proposed approach is stated to solve this problem in section 4.

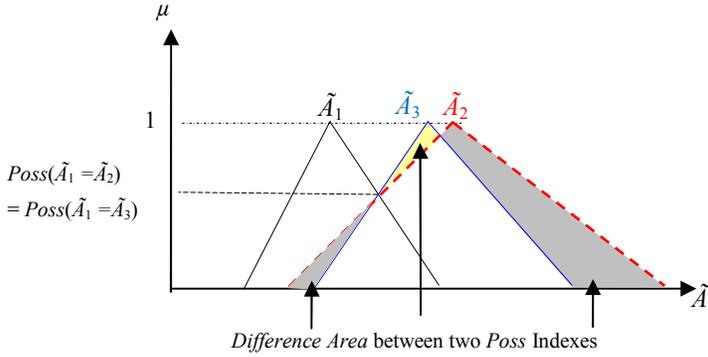


Figure 1: Equality Measure of three fuzzy numbers

This study introduces a new possibility equality index to establish new constraints in fuzzy regression analysis with an optimal confidence level, named h -level, in which conjunction problem is under consideration. Besides that a more precise method is described to calculate the amount of error calculation.

II- FUZZY LINEAR REGRESSION MODELS

Fuzzy Linear Regression (FLR) model was introduced initially by Tanaka et al. (1982) as:

$$\tilde{Y}_i^* = \tilde{A}_0 X_{i0} + \tilde{A}_1 X_{i1} + \dots + \tilde{A}_n X_{in} = \tilde{A} X_i \quad (3)$$

Where \tilde{Y}_i^* , $i=1, \dots, m$, are the estimated data, $\tilde{A}_j=(a_j, c_j)_L$, $j=0, 1, \dots, n$ are the set of symmetric fuzzy coefficients, and $X_i = [X_{i0}, X_{i1}, \dots, X_{in}]^T$ are the vector of independent variables.

The extension principle ([27]) plays basic role in the fuzzy set theory. It provides a fundamental for all manipulations on fuzzy sets. By applying it for the fuzzy linear regression model, the membership function of \tilde{Y}_i^* can be defined as:

$$\mu_{\tilde{Y}_i^*}(y^*) = \max_{y^*=f(X_i, \tilde{A})} \min_j (\mu_{\tilde{A}_j}(a_j)) \quad (4)$$

According to the extension principle, the optimization process is formulated as follows:

$$\min Z(h) = \sum_{i=1}^m \sum_{j=0}^n c_j |X_{ij}|$$

subject to

$$\sum_{j=0}^n a_j X_{ij} + |L^{-1}(h)| \sum_{j=0}^n c_j |X_{ij}| \geq y_i + |L^{-1}(h)| e_i, \quad i=1, 2, \dots, m \quad (5)$$

$$\sum_{j=0}^n a_j X_{ij} - |L^{-1}(h)| \sum_{j=0}^n c_j |X_{ij}| \leq y_i - |L^{-1}(h)| e_i, \quad i=1, 2, \dots, m$$

$$c_j \geq 0, a_j \in \mathbf{R} \quad j=0, 1, \dots, n.$$

Where $\tilde{Y}_i=(y_i, e_i)$, $i=1, 2, \dots, m$ are the fuzzy output observations. $|L^{-1}(h)|$ is supposed to be equal with $|L(h)|=1-h$; $0 < h < 1$, provided that the coefficients to be triangular fuzzy numbers. The role of h here is like a confidence level controller. In fact, it is close to zero in which more risk isn't acceptable and vice versa, it goes to one when the problem is considered at an optimistic point of view. Later, restriction, $c|X_i| \geq 0$; $i=1, \dots, m$, (where $c = [c_0, c_1, \dots, c_n]$) was substituted instead of $c \geq 0$ by (Change and Lee, 1994b), where m is the number of observations (Change and Lee [5]).

III- THE NEW APPROACH

A small change in constrictions regarding to new definition of the possibility of equality is taken into account here. Introducing new definition about the possibility equality index brings about some variation in constrains, especially the influence of h -level on the fuzzy regression model.

The conjunction problem, (6), guarantees that there will be always an overlap between the given outputs, \tilde{Y}_i , and the estimated fuzzy numbers, \tilde{Y}_i^* .

$$[\tilde{Y}_i^*]_h \cap [\tilde{Y}_i]_h \neq \emptyset \quad (6)$$

A possibility of equality definition is given in follow based on *Problem Definition* which was described above, to identify more precise estimated parameters with establishing new constraints in fuzzy linear regression models. For this reason the centers and spreads of two fuzzy numbers are regarded simultaneously.

Proposed Definition:

$$Poss(\tilde{A}_1 = \tilde{A}_2) = 1 - \frac{|a_2 - a_1| + |c_2 - c_1|}{|c_2 + c_1| + |c_2 - c_1|}$$

Where $\tilde{A}_1=(a_1, c_1)$ and $\tilde{A}_2=(a_2, c_2)$ are two symmetric triangular fuzzy numbers.

According to proposed definition, possibility of equality between two fuzzy numbers is one, if both of the following conditions are to be correct at the same time.

(i): $aX_i = y_i$

(ii): $c|X_i| = e_i$

Where aX_i and $c|X_i|$ are the center and spread of the estimated symmetric triangular fuzzy numbers, respectively ($\tilde{Y}_i = (aX_i, c|X_i|)$). The spreads of two fuzzy numbers are regarded as the criteria to compare them in which the first condition is right but the second is not. So that increasing of estimated spread, $c|X_i|$, with keeping constant of common area between two fuzzy numbers, causes to lessen the amount of possibility of equality, whereas this issue isn't seen in the possibility of equality index which has been introduced by (Dubois and Prade[7]). In other position, distance between the centers of two fuzzy numbers will be important factor when the spread of estimated data is to be equal with the spread of the output observation. Moreover, the maximum distance between aX_i and y_i (the centers of two fuzzy numbers) in which the *Conjunction* problem is applying, equals to $(c|X_i| + e_i)$. Thus, it guarantees that the numerator of fraction in proposed definition is less than its denominator and proposed index is ever less than or equal to one.

Lemma 1: with respect to proposed definition, $h_i = Poss(\tilde{Y}_i = \tilde{Y}_i^*) \geq h, i = 1, \dots, m$, if and only if

$$\sum_{j=0}^n a_i X_{ij} + (1-2|L^{-1}(h)|) \sum_{j=0}^n c_j |X_{ij}| \leq y_i + e_i, \quad i = 1, 2, \dots, m$$

$$\sum_{j=0}^n a_i X_{ij} - \sum_{j=0}^n c_j |X_{ij}| \leq y_i - (1-2|L^{-1}(h)|)e_i, \quad i = 1, 2, \dots, m$$

$$\sum_{j=0}^n a_i X_{ij} - \sum_{j=0}^n c_j |X_{ij}| (1-2|L^{-1}(h)|) \geq y_i - e_i, \quad i = 1, 2, \dots, m$$

$$\sum_{j=0}^n a_i X_{ij} + \sum_{j=0}^n c_j |X_{ij}| \geq y_i + (1-2|L^{-1}(h)|)e_i, \quad i = 1, 2, \dots, m$$

where $h = \min\{h_1, h_2, \dots, h_m\}$.

Proof:

$$poss(\tilde{Y}_i = \tilde{Y}_i^*) = 1 - \frac{|a^T X_i - y_i| + |c^T |X_i| - e_i|}{|c^T |X_i| + e_i| + |c^T |X_i| - e_i|} \geq h$$

$$1 - \frac{|a^T X_i - y_i| + |c^T |X_i| - e_i|}{|c^T |X_i| + e_i| + |c^T |X_i| - e_i|} \geq h$$

$$\frac{|a^T X_i - y_i| + |c^T |X_i| - e_i|}{|c^T |X_i| + e_i| + |c^T |X_i| - e_i|} \leq (1-h) = |L^{-1}(h)|$$

$$|a^T X_i - y_i| + |c^T |X_i| - e_i| \leq |L^{-1}(h)| \left\{ |c^T |X_i| + e_i| + |c^T |X_i| - e_i| \right\}$$

$$|a^T X_i - y_i| + |c^T |X_i| - e_i| (1 - |L^{-1}(h)|) \leq |L^{-1}(h)| |c^T |X_i| + e_i|$$

Last equation is divided into four cases which are listed as follows:

$$a^T X_i \geq y_i, \quad c^T |X_i| \geq e_i \Rightarrow a^T X_i + c^T |X_i| (1-2|L^{-1}(h)|) \leq y_i + e_i \quad (8)$$

$$a^T X_i \geq y_i, \quad c^T |X_i| < e_i \Rightarrow a^T X_i - c^T |X_i| \leq y_i - (1-2|L^{-1}(h)|)e_i \quad (9)$$

$$a^T X_i \leq y_i, \quad c^T |X_i| \geq e_i \Rightarrow a^T X_i - c^T |X_i| (1-2|L^{-1}(h)|) \geq y_i - e_i \quad (10)$$

$$a^T X_i \leq y_i, \quad c^T |X_i| < e_i \Rightarrow a^T X_i + c^T |X_i| \geq y_i + (1-2|L^{-1}(h)|)e_i \quad (11)$$

Objective Function: To achieve best objective function, it is necessary to consider the distance between centers, spreads and h -level, all together in order to get more reliable results. Shakouri G. and Nadimi [17] proposed an objective function based on Non-equality possibility index with considering the mentioned factors. It is applied here, which is given as follows:

$$\min \sum (|a X_i + L^{-1}(h) c |X_i| - y_i - L^{-1}(h) e_i| + |a X_i - L^{-1}(h) c |X_i| - y_i + L^{-1}(h) e_i|) \quad (12)$$

Based on (12) and Lemma 1, the optimization problem is summarized as follows:

$$\min \sum_{i=1}^m |a^T X_i + L^{-1}(h) c^T |X_i| - (y_i + L^{-1}(h) e_i)| + |a^T X_i - L^{-1}(h) c^T |X_i| - (y_i - L^{-1}(h) e_i)|$$

s.t.

$$\sum_{j=0}^n a_i X_{ij} + (1-2|L^{-1}(h)|) \sum_{j=0}^n c_j |X_{ij}| \leq y_i + e_i, \quad i = 1, 2, \dots, m$$

$$\sum_{j=0}^n a_i X_{ij} - \sum_{j=0}^n c_j |X_{ij}| \leq y_i - (1-2|L^{-1}(h)|)e_i, \quad i = 1, 2, \dots, m$$

$$\sum_{j=0}^n a_i X_{ij} - \sum_{j=0}^n c_j |X_{ij}| (1-2|L^{-1}(h)|) \geq y_i - e_i, \quad i = 1, 2, \dots, m$$

$$\sum_{j=0}^n a_i X_{ij} + \sum_{j=0}^n c_j |X_{ij}| \geq y_i + (1-2|L^{-1}(h)|)e_i, \quad i = 1, 2, \dots, m$$

$$h \leq 1,$$

$$c^T |X_i| \geq 0, \quad a, c \in R \text{ and } c \neq 0$$

Kim and Bishu [9]) proposed a criterion to evaluate fuzzy regression result, it is defined by the following index:

$$E_i = \frac{D}{\int_{S_{\tilde{Y}_i}} \tilde{Y}_i(y) \partial y}, \quad (14)$$

where D is difference between the two observed and estimated membership functions, which is obtained as follows:

$$D = \int_{S_{\tilde{Y}_i^*} \cup S_{\tilde{Y}_i}} |\tilde{Y}_i^*(y) - \tilde{Y}_i(y)| \partial y,$$

where $S_{\tilde{Y}_i^*}$ and $S_{\tilde{Y}_i}$ are the supports of observed value, \tilde{Y}_i , and estimated value, \tilde{Y}_i^* , respectively.

But there are some problems in this criterion, which are described as follows:

Denominator of (14) is an observed value which is constant. It can be estimated with different fuzzy numbers. For instance, Figure 2 shows three fuzzy numbers, that \tilde{A}_1 is an observed value, which has to be estimated. \tilde{A}_2 and \tilde{A}_3 can be dealt with as two options for the estimation of \tilde{A}_1 . As you can see, \tilde{A}_2 and \tilde{A}_3 have common areas with \tilde{A}_1 , but the spread of \tilde{A}_2 is wider than \tilde{A}_3 , while there is a little difference in the intersection area. According to mentioned criteria both estimated fuzzy numbers, \tilde{A}_2 and \tilde{A}_3 , may be assumed suitable for the observed fuzzy number, \tilde{A}_1 , with regard to some amount of error. Whereas to calculate the amount of error, not only it is required to consider the area of common region, but also it is necessary to take into account the unshared area, too. With regard to Kim and Bishu's criteria, the total amount of error for $\tilde{A}_1 \& \tilde{A}_2$ to $\tilde{A}_1 \& \tilde{A}_3$ isn't greater because there is no difference in their common area.

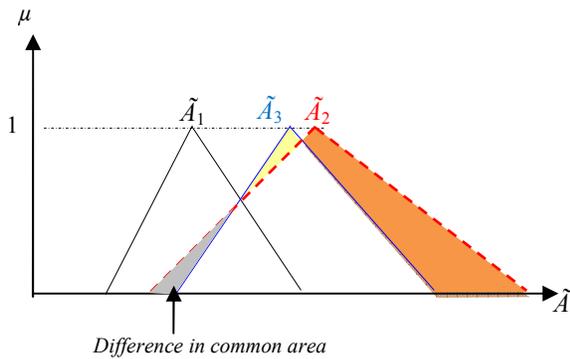


Figure 3: Compare Measure of three fuzzy numbers

To solve this problem which rises from the constant value of denominator in E_1 , it is necessary to consider estimated and observed membership functions, together in the denominator of E_1 .

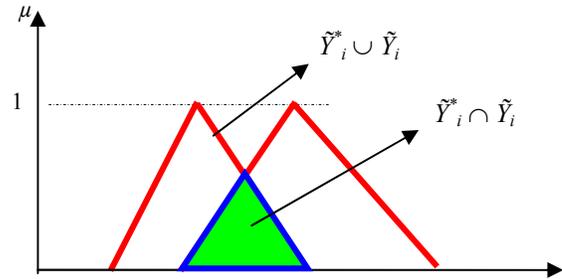


Figure 4: Concept of Union and Intersection in fuzzy numbers

So union and intersection concept (as shown in Figure 4) is applied to compare and evaluate two fuzzy numbers equality as follows: (for more information, see [26], [27])

$$E_i = 1 - \frac{\min \{ \mu_{\tilde{Y}_i^*}(y), \mu_{\tilde{Y}_i}(y) \}}{\max \{ \mu_{\tilde{Y}_i^*}(y), \mu_{\tilde{Y}_i}(y) \}} \quad (16)$$

By adding the union concept in the denominator of (16) the abovementioned problem can be solved. The numerical examples will confirm the improvement of the model.

IV- NUMERICAL EXAMPLE

Four examples are considered to evaluate the new approach and to compare them with other methods. Tanaka's method (TM) (Tanaka et al.[24]), Savic and Pedrycz method (SP) (Savic and Pedrycz [15]), Kim and Bishu method (KB) (Kim and Bishu [9]), Modarres Approach (MA) (Modarres et al. [12]), Shakouri and Nadimi method (NA) (Shakouri G. and Nadimi[17]) are chosen to evaluate and compare with the New Approach (*NAI*).

Example 1. Tanaka's data (Tanaka et al.[24]), shown in Table 1, is used as the first example. Applying *NAI* to these data, the following fuzzy regression model is obtained:

$$\tilde{Y}^* = (5.209, 4.962)_L + (1.558, -0.512)_L x, \quad h = 0.535$$

where the fuzzy coefficients are given according to the format of $\tilde{A}_j = (a_j, c_j)_L$. Corresponding errors of this method is also given in the same table, where it is compared to the other methods. It is demonstrated that the total error of *NAI* is 2.469, which is the least as compared to all other methods, which shows better performance.

Table 1: Original Data and the Estimation Errors for Example 1

<i>I</i>	x_i	(y_i, e_i)	Errors in estimation					
			TM	SP	KB	MA	NA	<u>NAI</u>
1	1	$(8,1.8)_L$	0.745	0.691	0.735	0.780	0.746	0.699
2	2	$(6.4,2.2)_L$	0.640	0.687	0.827	0.780	0.804	0.692
3	3	$(9.5,2.6)_L$	0.431	0.467	0.369	0.205	0.259	0.322
4	4	$(13.5,2.6)_L$	0.519	0.605	0.729	0.782	0.804	0.756
5	5	$(13.0,2.4)_L$	0.569	0.512	0.328	0.091	0.000	0.000
Total errors			2.905	2.962	2.988	2.962	2.614	2.469

Table 2: Data for Example 2

	Response time	Inside control room experience	Outside control room experience	Education
Team 1	$(5.83, 3.56)_L$	2.0	0.0	15.25
Team 2	$(0.85, 0.52)_L$	0.0	5.0	14.13
Team 3	$(13.93, 8.5)_L$	1.13	1.5	14.13
Team 4	$(4, 2.44)_L$	2.0	1.25	13.63
Team 5	$(1.65, 1.01)_L$	2.19	3.75	14.75
Team 6	$(1.58, 0.96)_L$	0.25	3.5	13.75
Team 7	$(8.18, 4.99)_L$	0.75	5.25	15.25
Team 8	$(1.85, 1.13)_L$	4.25	2.0	13.5

Table 3: Comparison between Estimation Errors for Example 2

Team Number	Errors in estimation					
	TM	SP	KB	MA	NA	<u>NAI</u>
Team 1	0.84	0.82	0.76	0.87	0.58	0.00
Team 2	0.99	0.99	1.00	0.79	0.34	0.19
Team 3	0.26	0.57	0.92	0.91	1.00	1.00
Team 4	0.92	0.81	0.44	0.43	0.36	0.33
Team 5	0.93	0.95	0.92	0.83	0.93	0.55
Team 6	0.98	0.96	0.97	0.83	0.03	0.17
Team 7	0.60	0.63	0.84	0.92	1.00	1.00
Team 8	0.93	0.87	0.72	0.75	0.00	0.00
Total errors	6.45	6.61	6.56	6.34	4.25	3.23

Example 2. Kim and Bihu (Kim and Bishu [9]) applied data relative to the nuclear power plant control room crew. The

following fuzzy regression model is obtained by applying NAI approach to the data given in Table 2:

$$\tilde{Y}^* = - (10.499, 9.425)_L - (0.192, 0.399)_L x_1 - (0.816, 0.421)_L x_2 + (1.096, 0.857)_L x_3$$

Which $h = 0.0010$ is the optimum value. The error for each sample output and the total error are obtained by (16) and illustrated in Table 3.

The results show that the proposed method, NAI, has considerably reduced the total error in comparison with other methods.

Example 3. Shakouri and Nadimi [17] considered an example to compare their approach with other methods based on the fuzzy regression for which the fuzzy parameters were available, as well as the inputs. The example is given as:

$$\tilde{Y} = (2., 1)_L + (3, 0.5)_L x \quad (17)$$

where:

$$A_0 = (a_0, c_0)_L = (2, 1)_L, \quad A_1 = (a_1, c_1)_L = (3, 0.5)_L$$

Therein the fuzzy regression model is regarded in diverse way. In brief, they presupposed that fuzzy regression model is available and then they estimated the parameters of the model to calculate error terms and accuracy of their approach. Here the results of proposed approach are compared with other methods. The fuzzy outputs, however, can be calculated by the given model. NAI approach is used to compare with the NA and MA methods. Results for $\tilde{Y} = (y, e)_L$ are given in Table 4 based on corresponding x .

Table 4: Crisp inputs and the corresponding fuzzy outputs for model (17)

X	y	E
1	5	1.5
2	8	2
3	11	2.5
4	14	3
5	17	3.5

Three methods are applied to estimate the fuzzy parameters and the results of each approach are summarized in the following table:

Table 5 : Fuzzy Parameters Estimates by MA, NA and NAI

Method	$A_0 = (a_0, c_0)_L$	$A_1 = (a_1, c_1)_L$
MA	(2.000000, 1.000000) _L	(3.000000, 0.500000) _L
Parameters estimated by MA	(1.998005, 1.840016) _L	(3.000000, 0.160000) _L
Parameters estimated by NA	(2.000000, 1.000111) _L	(3.000000, 0.500089) _L
Parameters estimated by <u>NAI</u>	(2.000000, 1.000000) _L	(3.000000, 0.500000) _L

Original Parameters	(2.000000, 1.000000) _L	(3.000000, 0.500000) _L
Parameters estimated by MA	(1.998005, 1.840016) _L	(3.000000, 0.160000) _L
Parameters estimated by NA	(2.000000, 1.000111) _L	(3.000000, 0.500089) _L
Parameters estimated by <u>NAI</u>	(2.000000, 1.000000) _L	(3.000000, 0.500000) _L

Here, MA method is calculated with a tolerance of $\varepsilon = 0.000001$ to find the optimal h -level by iteration Modarres et al. (2005). In this example the amount of h for MA, NA and NAI, is equal to 0.9985722, 0.8623157 and 0.2512473, respectively. Table 5 shows that the proposed approach is closer to the real data as compared to MA. There are however a few differences between NA and NAI approaches. C_0 and c_1 are two cases with differences of 0.000111 and 0.000089, respectively. Based on the results which were demonstrated in Table 5, NAI approach is more precise than NA method.

Example 4. In this example MA method is left out, NA and NAI results are taken into account and compared as follows. Tanaka's data, (Tanaka [22]) which has been shown in Table 6, is considered here to evaluate and compare the NAI approach with NA method.

Table 6: crisp input and fuzzy output with fuzzy error

x_1	x_2	x_3	Y	e
3	5	9	96	42
14	8	3	120	47
7	1	4	52	33
11	7	3	106	45
7	12	15	189	79
8	15	10	194	65
3	9	6	107	42
12	15	11	216	78
10	5	8	108	52
9	7	4	103	44

Herein, at first, fuzzy linear regression model is considered with the assumption of \tilde{A}_0 , afterwards it is dropped off the model. Equations (18) and (19) demonstrate the results of NA and NAI approaches, respectively.

$$\tilde{Y} = (2.9574, 2.8893)_L x_1 + (8.2470, 1.2419)_L x_2 + (4.6636, 2.3257)_L x_3 \quad (18)$$

$$\underline{NAI}, \quad \tilde{Y} = (2.9712, 2.8793)_L x_1 + (8.2057, 1.0338)_L x_2 + (4.7143, 2.6460)_L x_3 \quad (19)$$

Above fuzzy linear regressions were achieved without considering \tilde{A}_0 . The model parameters are estimated once more by the NA and NAI methods with inclusion of \tilde{A}_0 which are given as follows:

$$NA \quad \tilde{Y} = (8.2670, 0.4792)_L + (2.3613, 2.0617)_L x_1 + (8.1775, 0.9281)_L x_2 + (4.4179, 3.4107)_L x_3, \quad (20)$$

$$\underline{NAI} \quad \tilde{Y} = (7.1071, 13.3249)_L + (2.4299, 1.2944)_L x_1 + (8.1821, 0.9020)_L x_2 + (4.4722, 2.7789)_L x_3, \quad (21)$$

Table 7 demonstrates the error of each approach in two different states. It is evident from the data that there are some differences between NA and NAI approaches. Here,

the total errors for the proposed approach in each state are also less than the NA method.

Table 7 : Comparison between Estimation Errors with two methods

Number	Errors in estimation			
	With \tilde{A}_0		Without \tilde{A}_0	
	<u>NAI</u>	NA	<u>NAI</u>	NA
1	0.0000	0.1012	0.1043	0.1474
2	0.0000	0.0000	0.1886	0.2000
3	0.1426	0.1008	0.2202	0.2350
4	0.1136	0.0619	0.0394	0.0517
5	0.0251	0.0566	0.0968	0.1207
6	0.0000	0.0000	0.0000	0.0000
7	0.1552	0.0174	0.1948	0.1953
8	0.1859	0.1831	0.1265	0.1290
9	0.0193	0.0191	0.0568	0.0318
10	0.1097	0.0363	0.0061	0.0000
Sum of the Output Errors	0.57624	0.75137	1.03338	1.11089

V- CONCLUSION REMARKS

Possibility of equality was developed with its application in fuzzy regression analysis in this study. It was then used to formulate new constraints in fuzzy regression model based on conjunction problem. Meanwhile an accurate criterion, based on union and intersection concept, was introduced to assess a fuzzy regression result. The results of examples confirmed the accuracy and improvement of proposed approach as compared to the other methods.

VI- REFERENCE

- [1]. Al-Kandaria, A.M., S.A. Solimanb, and M.E. El-Hawary. "Fuzzy short-term electric load forecasting." *Electrical Power and Energy Systems* 26, 2004: 111-122.
- [2]. Beenstock, M., E. Goldin, and D. Nabot. "The demand for electricity in Israel." *Energy Economics* 21, 1999: 168-183.
- [3]. Carcedo, J.M., and J.V. Otero. "Modeling the non-linear response of spanish electricity demand to temperature variation ." *Energy Economics* 27, 2005: 477-494.
- [4]. Change, P.T., and E.S. Lee. "Fuzzy least absolute deviation regression and the conflicting trends in fuzzy parameters." *Computational Mathematical Application* 28 (5), 1994: 89-101.
- [5]. Change, P.T., and E.S. Lee. "Fuzzy linear regression with spreads unrestricted in sign." *Computational Mathematic Application* 28 (4), 1994: 51-70.
- [6]. Diamond, P. "Fuzzy least squares." *Information Science* 46, 1988: 141-157.
- [7]. Dubois, D., and H. Prade. *Fuzzy Sets and Systems: Theory and Applications*. New York: Academic Press, 1980.
- [8]. D'Urso, P., and T. Gastaldi. "A least-squares approach to fuzzy linear regression analysis." *Computational Statistics & Data Analysis* 34, 2000: 427-440.
- [9]. Kim, B., and R.R. Bishu. "Evaluation of fuzzy linear regression models by comparison membership functions." *Fuzzy Sets and Systems* 100, 1998: 343-352.
- [10]. Manusov, V.Z., and A.V. Mogilenko. "The fuzzy regression analysis as a means of electric power losses evaluation in electrical networks." *Power System Management and Control. IEEE*, 2002: No. 488 .
- [11]. Ming, M., M. Friedman, and A. Kandel. "General fuzzy least squares." *Fuzzy Sets and Systems* 88, 1997: 107-118.
- [12]. Modarres, M., E. Nasrabadi, and M.M. Nasrabadi. "Fuzzy linear regression models with least square errors." *Applied Mathematics and Computation* 163, 2005: 977-989.
- [13]. Nazarko, J., and W. Zalewski. "The Fuzzy Regression Approach to Peak Load Estimation in Power Distribution Systems ." *Transactions on Power Systems. IEEE*, 1999: No. 3.
- [14]. Olsina, A.F., F. Garces, and H.J. Haubrich. "Modeling long-term dynamics of electricity markets." *Energy Policy* 34, 2006: 1411-1433.
- [15]. Savic, D.A., and W. Pedrycz. "Evaluation of fuzzy linear regression models." *Fuzzy Sets and Systems* 39, 1991: 51-63.
- [16]. Sedelnikov, A.V., and V.Z. Manusov. "The estimate of loss of electric power under conditions of uncertainty by using fuzzy regression analysis with asymmetrical triangular coefficients." *Electrical Engineering. IEEE*, 2004.
- [17]. Shakouri G., H., and R. Nadimi. "A novel fuzzy linear regression model based on a non-equality possibility index and optimum uncertainty." *Applied Soft Computing* 9, 2009: 590-598.
- [18]. Shakouri G., H., and R. Nadimi. "Investigation on the short-term variations of Electricity Demand due to the Climate Changes via a Hybrid TSK-FR Model." *Industrial Engineering & Engineering Management. Singapore: IEEE*, 2007: 807-811.
- [19]. Shakouri G., H., J. Nazarzadeh, and S.K.Y. Nikravesh. "Exogeneity Investigation and Modeling Energy Demand via Parallel Dynamic Linear Models for Maximum Simultaneous Power Demand." *Conference on Control Applications. IEEE*, 2003.
- [20]. Shakouri G., H., M. Rastad, and J. Nazarzadeh. "Exogeneity A Hybrid Nonlinear Model for the Annual Maximum Simultaneous Electric Power Demand." *IEEE* 21 2006: No.3.
- [21]. Shakouri G., H., R. Nadimi, and F. Ghaderi. "A hybrid TSK-FR model to study short-term variations of the electricity demand versus the temperature changes." *Expert Systems with Application* 36, 2009: 1765-1772.
- [22]. Tanaka, H. "Fuzzy data analysis by possibilistic linear models." *Fuzzy Sets and Systems* 24, 1987: 363-375.
- [23]. Tanaka, H., and J. Watada. "Possibilistic linear systems and their application to the linear regression model." *Fuzzy Sets and Systems* 27, 1988: 275-289.
- [24]. Tanaka, H., I. Hayashi, and J. Watada. "Possibilistic linear regression analysis for fuzzy data." *European Journal of Operational Research*, 1989: 389-396.
- [25]. Tanaka, H., S. Uegima, and K. Asia. "Linear regression analysis with fuzzy model." *Transactions Systems Man Cybernet* 12, 1982: 903-907.
- [26]. Zadeh, L.A. "Fuzzy sets as a basis for a theory of possibility." *Fuzzy Sets and Systems* 1, 1978.
- [27]. Zimmermann, H.J. In *Fuzzy Set Theory and Its Application*, 33-35. Kluwer Academic Publishers, 1996.

Strategies and methodologies of Experimental Design in the online environment

by

Teresa Oliveira and Amilcar Oliveira

CEAUL and DCeT - Universidade Aberta, Lisboa

Abstract

Experimental Design is a branch of research in many areas, very varied structures and with many answers not yet known, being one of the most fascinating fields of research in Statistics. It has underlying ideas as important and in vogue as the optimization of factors, models and features, quality and competitiveness. It is a current powerful technique, indispensable in any experience, either in the definition of data to study – what type of data and how much data, or to choose the method and conditions of gathering the samples, always looking for the maximization of feedback information and minimizing costs.

Experimental Design applications are known from experiments in areas as diverse as Medicine, Engineering, Cryptography, Bioinformatics, Social Sciences and Education Sciences.

The technological innovations of today allowed prodigious advances in all areas of research and in particular at the level of Statistics and Experimental Design. Besides the usual computer programs such as STATISTICA, SPSS and SAS, with a relevant role in the programs of classroom teaching, researchers and teachers felt the need to create simple software, free, open to the community and manageable according to the specific needs in each case. R emerges as the current program for more investment in the scientific community of Statistics, making it especially attractive in education programs online. This paper investigates strategies and methods of Experimental Design, as well as R developments, aimed at applications in e-Learning/e-Teaching of these important themes in Masters courses online. Examples of experiences at UAb-Portugal will be presented.

1. Introduction

The importance of Statistics and Experimental Design in the education of future professionals in the various fields of scientific and technological knowledge is very well known as a current challenge. In general, teaching in these areas tends to be increasingly based on the use of Web resources and Software for specific assistance, because they are both popular and quite attractive, and also because they allow the development of an experimental and interacting components to the learning process, that was largely absent from the pedagogical tools and approaches previously used. We observe the online teaching programs with a growing trend, so it urges to pay special attention to developments aiming at new features which may complement the more traditional way of teaching. To this end, we will review the main resources available to support the online teaching of Experimental Design topics, emphasizing the role of Software R in this context. We will try to point out ways that lead to good practices for the future in this area, presenting the brief history of e-learning/e-teaching in the Master Course on Statistics, Maths and Computation at the Universidade Aberta (UAb), Portugal.

2. The impact of current Experimental Design

The strong impact of Experimental Design currently is mainly due to the fact that it has underlying ideas as important and in vogue as quality, competitiveness and optimization of resources, of factors and of models. It is a powerful methodology, indispensable in any experience, either in the definition of samples and data to study, what type of data and sample size, whether in choosing the method and conditions of sampling. The main objectives of Experimental Design are the maximization of information response and the minimization of the costs involved. The role on well designed experiences is crucial and applications of Experimental Design are known in areas as diverse as Medicine, Engineering, Cryptography, Bioinformatics, Social Sciences and Education Sciences.

In the twenty-first century we have been assisting to the fact that many of the best-qualified jobs place as a reality the need of developing new skills in mathematical, statistical and computational abilities. So, whether in developed or developing countries, we assist in the last decade to a trend of increasing demand in respect of courses and subjects related to Statistics in general and Experimental Design in particular.

Simultaneously, in the last decade we are witnessing the success of the impact of online learning, emerging to address the difficulties inherent in traditional presence and in traditional distance education, now allowing students to follow lectures in any place, at anytime and at a very affordable cost.

All these developments and changes are very recent, and although there are many strong research teams looking for new results and advances in Experimental Design, less attention has been paid to the need of developing new tools and new methodologies to better accomplish the teaching of these techniques, according to the new century trends and demands of E-learning/E-teaching courses. This is a very challenging new field of research with many open problems to be known and answered.

It was already noticed besides from short time experience that, for students on the online learning of statistics and experimental design topics, it is very stimulating to develop interactive activities, involving student-student and student-teacher-student interactions. Particularly in Master courses the role of the teacher increasingly assumes the character of a companion study, of course in parallel with the transmission of knowledge, but always encouraging the process of self-learning. The role of software R becomes here fundamental. Software R proves to be a powerful tool, bridging the needs felt by researchers and teachers on the creation of free software, open to the community, simple and manageable according to the specific needs in each case.

3. E-learning/E-teaching Experimental Design: Tools and computational resources

Currently Experimental Design is one of the most fascinating fields of research, providing a powerful challenge and opportunity to obtain new results on theses and on research projects with links to many different areas. This leads to the increasing development of Experimental Design issues in Graduate and Master Courses, not only in Maths and Statistics, but also in areas such as Engineering, Environmental Sciences, Health Sciences and Feeding Consuming Sciences.

The use of Web and of computational resources to support education deserves special attention, especially when it comes to online learning, since in this scenario it will certainly be an important complement to the student. We will present a retrospective summary of the main resources currently available on the Web to support the online teaching of Experimental Design, especially those of free access. Currently it is easy to

find online interesting virtual labs to support teachers and students, interactive applets as tools of great potential for learning, electronic books and the particular software adequate to the e-learning/e-teaching of Statistics and Experimental Design - R: a free open source for Statistics learning.

Among other available e-books in Statistics, MD*BOOKS site present a very good selection on statistical subjects at the website:

<http://www.xplore-stat.de/ebooks/ebooks.html> .

Canavos and Koutrouvelis(2009) present an e-book on the introduction to the Design and Analysis of Experiments and some other important e-books on Experimental Design can be found at the websites:

<http://www.math.vu.nl/sto/onderwijs/doeanova/notes.pdf>;

<http://instructors.coursesmart.com/0136158633> ;

<http://www.itl.nist.gov/div898/handbook/index.htm>;

<http://www.itl.nist.gov/div898/handbook/pri/pri.htm>

Particular emphasis to the potential of R software on the support to the teaching of Experimental Design in the online environment is now stated together with a brief introduction. Information about the project and how to download the program, as well as sources of documentation are available at <http://www.r-project.org> .

R is an integrated project involving, among other means, a language and environment for statistical computing. As a part of the GNU Project it is an affordability free source to the entire scientific community, and the community that supports it, provide access to state-of-the-art of statistical graphics, visualization and computing, considering many levels of users expertise. This project was initiated and developed from the language S by Ross Ihaka and Robert Gentleman from the 1990s. R provides a wide range of resources, from packages developed by researchers on a worldwide network, with applications in most areas of science, e-books in several languages and on different themes in the field of Statistics and Experimental Design. It presents an integrated environment for data manipulation, calculations and graphical representations and is highly extensible.

Advantages of using R on online teaching programs, instead of classical computer programs such as STATISTICA, SPSS and SAS are then obvious, mainly considering the expensive cost licences, difficulting their usage at public education institutions.

In what concerns to Experimental Design issues using R the website <http://www.stat.washington.edu/fritz/DATAFILES/Stat421Rintro.pdf> presents a brief and very simple introduction and at http://cran.r-project.org/doc/contrib/Vikneswaran-ED_companion.pdf we find the basic methodologies of Experimental Design using R, which gives an important help on the students approach.

R has a nice amount of functionality for Experimental Design or Design of Experiments (DOE) which appears in various R packages. Gromping, U. (2008-2009) present the CRAN Task View on Design of Experiments, available at <http://cran.r-project.org/web/views/ExperimentalDesign.html>. Several packages were developed concerning to solve DOE problems, such as AlgDesign which creates full Factorial Designs and Mixture Designs, among others; Conf.design which is a package adequate to create a design with certain interaction effects confounded with blocks, allowing combine designs in several ways; Package blockTools which is adequate to assign units to blocks in a Block Design and Package agricolae which was especially developed to solve agricultural and plant breeding experiments. Some further packages handle special situations in DOE, but still some Designs have open fields to look for adequate software developments.

Baier and Neuwirth (2007) refer to the convenience of integrating R in Microsoft Excel since this provides a good way to combine the advantages of a spreadsheet with the flexibility of R. This seems very interesting since it will help to provide solutions for some non-regular situations of Designs for which still there are not yet R convenient available packages.

4. Practices and strategies on teaching Experimental Design topics online: examples

We present strategies adopted to teach Experimental Design online, considering classes with a big number of students and classes with not so big number of students, which

sometimes allow to go a bit further and fit students needs according to their professional fields.

In the first case, we refer to Darius and Schrevens (2006), where a very interesting experience was presented. The authors alert to that students typically have little opportunity to get experience in the ability to design experiments, since in most of the courses there is more emphasis on the analysis of data already collected than on the actual design process. Also usually in classes, exercises are presented to students with data supplied, and lack of information on the reasons that led to the particular way of obtaining data. This concern is also present in several classical books, as Dean and Voss (1998) and Montgomery (2009). In literature many authors stressed the importance of including projects into the courses, in which students have to perform and analyse a real experiment, but such projects are unfeasible for classes with a big number of students. To go over this situation Darius and Schrevens (2006) present as strategy the use of virtual experiments in teaching design and analysis of experiments. A tool is explored for gaining design experience: computer virtual experiments. This consist in “software environments which mimic a real situation of interest, pose a research question, then invite the user to collect associated data which, when statistically analysed, will shed light on the research question”. A collection of sorts referred to as ENV2EXP, available at <http://www.kuleuven.be/ucs/env2exp> have been experimented and three of the applets were discussed in this work: The Factory Applet (adequate to study designs such as Fractional Factorial or other Screening Designs, Full Factorial Designs, Box-Behnken) ; the Greenhouse Applet (this applet allows the user to get comparative experience with almost all classical designs, as completely randomized, complete or incomplete block designs and Latin Square Designs) and the Shooting Applets (very simple, can be used in the context of an introductory statistics course). With these applets is was shown that nowadays technology allows the creation of accessible and rich environments, adequate to improve the online teaching, providing students with new experience skills. Darius and Schrevens (2006) recommend a companion collection of programs referred to as VESTAC (Darius et al., 2000, and at <http://www.kuleuven.be/ucs/java>) to illustrate many associated statistical concepts.

In the second case we present and discuss our work experience on teaching Experimental Design issues at Universidade Aberta - UAb, Portugal in a Master course for the last three

years. After a long period of teaching in distance education in its most traditional form, the UAb implemented a new pedagogical online teaching model in 2007, which uses Moodle - Modular Object Oriented Dynamic Learning Environment- as a platform to teach and conduct under graduate and Master courses online. The yet short experience so far already allows us to draw some conclusions and ideas for the future. The Department of Science and Technology is currently responsible for teaching the third Edition of the online Master in Statistics, Maths and Computation (MSMC). One of this Masters Course options is Computational Statistics and for this we have the following curricula:

Table 1 Master in Statistics, Maths and Computation - Computational Statistics: A structure overview

Semester I	Semester II
Statistics I	Statistics II
Sampling and Data Analysis	Multivariate Data Analysis and Applications
Quality Control	Numerical Methods
Statistical Computation I	Statistical Computation II
	Significant Learning of Sciences

On the MSMC – Computational Statistics, semester II, in Statistics II some topics of Experimental Design are presented, based in using R packages to solve practical problems.

To illustrate R simplicity we present a generic exercise of a Randomized Block Design in which R and package STAT was used on variety comparison with ANOVA.

Table 2: Results on melon yield (Kg by plot), in an experiment to compare 4 varieties (Oliveira, T.A. 2004, pg. 279)

Variety	Blocks		
	I	II	III
Green Plain	130.0	130.2	131.5
Green Rough	131.8	125.3	131.2
White Plain	138.4	146.7	145.7
White Rough	131.3	130.0	129.7

Instructions in R should follow the input:

```
> y<-scan()
```

```
1: 130.0
2: 130.2
3: 131.5
...
10: 131.3
11: 130.0
12: 129.7
13:
Read 12 items
```

Then it follows the construction of a *data.frame*, with the data and indicators for blocks and varieties.

```
exel<-data.frame(variedade=factor(rep(1:4, each=3)), bloco=factor(rep(1:3,
4)), resp=y)
```

In a first approach exploring data, we have the Box-Plot:

```
> names(exel)
[1] "variedade" "bloco"      "resp"

> summary(exel)
variedade bloco      resp
1:3        1:4  Min.    :125.3
2:3        2:4  1st Qu.:130.0
3:3        3:4  Median :131.2
4:3                Mean  :133.5
                3rd Qu.:133.4

> attach(exel)
> plot(resp~variedade+bloco)
```

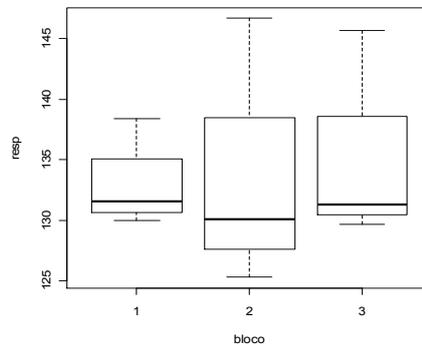


Figure 1: Box Plot for the Melon yield example

And to obtain the ANOVA Table, R instructions follow:

```
> exel.av<-aov(resp~bloco+variedade)
> anova(exel.av)
Analysis of Variance Table

Response: resp
      Df Sum Sq Mean Sq F value Pr(>F)
bloco  2   6.57    3.29  0.3126 0.74278
variedade 3 411.54  137.18  13.0505 0.00487 **
Residuals 6  63.07   10.51
```

Students are encouraged to undertake special projects to gain experience in design, data analysis and statistical computing. In this Master we have been observing a growth

trend in the number of students: first edition with 8 students, second edition with 20 students and third edition with 33 students. However in the optional class of Statistics II, second semester, still the number of students in the third edition is not too big (8), which still allow us to adopt some strategies not feasible otherwise, namely in what concerns to the last activity proposed to the students. Along the online classes students are introduced to the learning activities, in which they are invited to actively participate in Forums and to do some collaborative and web research on Experimental Design issues, according to the program: Introduction to Experimental Design; Fixed, Random and Mix Models in Experimentation; Complete and Incomplete Block Designs; Factorial Designs; Fractional Factorial Designs; Response Surface Methodologies and Advanced Experimental Design Models. Exercises are proposed and through the resolution of problems and reflexion supported by experimentation, available computational resources are explored, namely in what concerns using R and the DOE package. Some activities are individual and some are to be solved in group, as it is explained to the students in a previous schedule of activities, currently described in the platform webpage. The Moodle system, as screenshots is depicted in Figure 2.

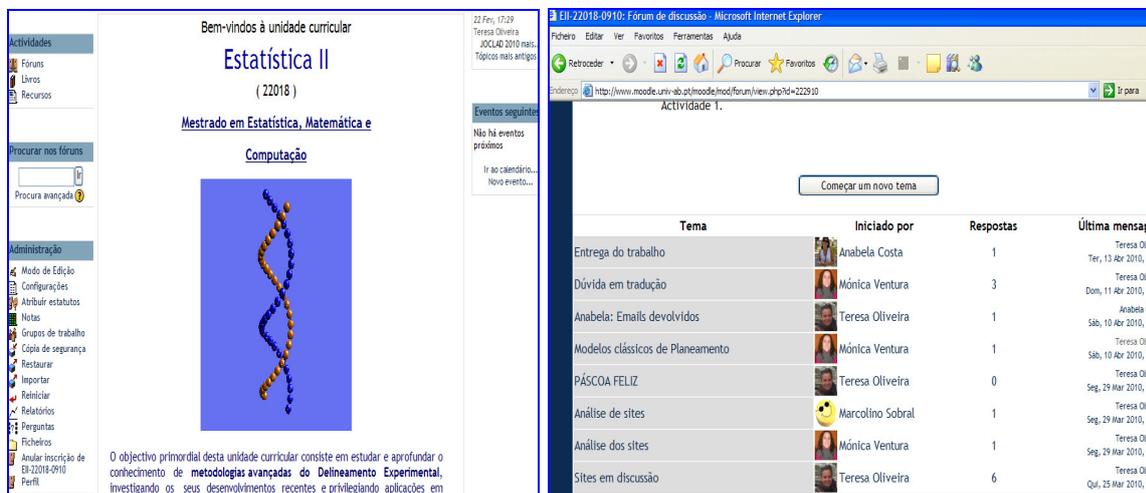


Figure 2 Environments in the Moodle system

The methodological strategies applied during the activities development include the use of real data bases, the resolution of problems presented in the adopted books by using R, as well as simulations and visualizations. In the last activity we promote collaboration and reflection on students own ideas and experiences to stimulate the use of the acquired

knowledge as a tool to improve their professional development and skills. Students are asked not only to perform and analyse a real experiment - using at least one of the DOE tools explored, but are also invited to suggest an application according to their professional fields. With the use of this strategy we provide a learning environment in which we can be assured that students will better understand the problem and surely will be motivated to solve it. As very interesting results of applying this strategy for three years, some master thesis projects were designed, some papers were presented by the students in national meetings and even some collaboration projects between the University and some Public Institutions are going on.

5. Considerations and Perspective Research

In a highly competitive and evolving world it seems crucial fostering the interest and involvement of research teams and universities on the solution of real world problems experienced in various relevant public and private institutions, by designing surveys and experimental designs, by developing and investigating stochastic models and computer simulations. The adoption of teaching strategies aimed at solving real problems will stimulate the professionals from many areas to look for Universities with the aim of improving their skills.

In a prospective research it's our aim to develop an experimental design in order to identify significant differences between using R and using other classical software for statistical simulations on student's knowledge and skills achievement, considering the two SMC-Master areas at UAb: Computational Statistics and Computational Mathematics.

References

- Baier, T. and Neuwirth, E.(2007). Excel::COM::R. *Computational Statistics*, 22(1), 91-108.
- Canavos, G.C. and Koutrouvelis, J.A.(2009). Introduction to the Design and Analysis of Experiments, Prentice Hall.
- Crawley, M. J. (2007). The R Book. John Wiley & Sons-New York. ISBN-13: 978-0-470-51024-7.
- Darius, P. and Schrevens, E. (2006). Use of Virtual Experiments in teaching Design and Analysis of Experiments. Proceedings of ICOTS-7.
- Darius, P.L., Ottoy, J-P., Solomin, A., Thas, O., Raeymaekers, B., and Michiels, S. (2000). A collection of applets for visualizing statistical concepts. In J.G. Bethlehem and

P.G.M. van der Heijden (EDS.), Proceedings in Computational Statistics 2000, Utrecht, The Netherlands. Heidelberg: Physica Verlag.

Dean, A. and Voss, D. (1998). Design and Analysis of Experiments. New York: Springer.

Mills, J.D. (2002). *Using Computer Simulation Methods to Teach Statistics: A review of the literature*, Journal of Statistics Education, 10(1).

Montgomery, D. (2009). Design and Analysis of Experiments, 7th Edition. Wiley.

Oliveira, T.A.(2004) . Estatística Aplicada. Edições Universidade Aberta, Portugal.

Oliveira, A. & Oliveira, T.A. (2008). *E-Learning and E-Teaching Statistics: Moodle and R Applications*. 11-14 June 2008, Lisbon. Article publishing in CD, EDEN 2008 ANNUAL CONFERENCE- New Learning Cultures, Edited by Alan Trait and András Szucs on the behalf of the European Distance and E-Learning Network.

R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, (2007). URL <http://www.r-project.org/> . ISBN 3-900051-07-0.

WebSites

1. <http://www.xplore-stat.de/ebooks/ebooks.html>, (*Research Data Center*), cited 8 March 2010
2. <http://mo161.soci.ous.ac.jp/@d/index.html>, (*Data oriented Statistical System*), cited 10 March 2010
3. <http://www.statistiklabor.de/en/index.html>, (*The Statistical Lab*), cited 10 March 2010
4. <http://onlinestatbook.com/rvls/index.html>, (*Rice Virtual Lab in Statistics*), cited 10 March 2010
5. <http://www.stat.tamu.edu/~jhardin/applets/index.html>, (*Java Applets*), cited 10 March 2010
6. <http://surfstat.anu.edu.au/surfstat-home/surfstat-main.html>, (*SurfStat Australia*), cited 10 March 2010
7. <http://wise.cgu.edu/index.html>, (*WISE Web Interface for Statistics Education*), cited 8 March 2010
8. <http://www.amstat.org/publications/jse/>, (*Journal of Statistics Education*), cited 8 March 2010
9. <http://www.educyclopedia.be/education/mathematicsjavastat.htm>, (*Educyclopedia, The Educational Encyclopedia*), cited 11 March 2010
10. <http://www.mhsatman.com/applets.php>, cited 11 March 2010
11. <http://statlink.tripod.com/>, (*Statlink*), cited 11 March 2010
12. <http://office.microsoft.com/pt-pt/excel/default.aspx>, cited 5 March 2010
13. <http://www.statsoft.com/>, cited 10 March 2010
14. <http://www.moodle.org>, cited 5 March 2010
15. <http://www.univ-ab.pt>, cited 30 April 2010
16. <http://r-project.org>, cited 3 May 2010
17. <http://lstat.kuleuven.be/java/> (*Java Applets for Visualization of Statistical Concepts*), cited 10 March 2010

18. <http://www.math.uah.edu/stat/> (*Virtual Laboratories in Probability and Statistics*), cited 10 March 2010
19. <http://socr.ucla.edu/> (*Statistics Online Computational Resources*), cited 3 May 2010

Test for the Absence of Clusters in 1-Dimensional Observations

Leonid Morozensky¹, Victor Olman², Yanbin Yin², Zeev Volkovich¹,
and Ying Xu²

¹ORT Braude Academic College of Engineering, Karmiel, Israel

Email: leon407@013net.net

²Department of Biochemistry and Molecular Biology, Computational Systems Biology Lab (CSBL), University of Georgia, Athens, USA

Email: olman@bmb.uga.edu

Abstract: A new statistical test for uniformity based on the sum of squared distances between neighboring one-dimensional observations is proposed. Analytical results for the statistics moments are presented as well as a computational recursive procedure for p-value calculations. Examples of testing the uniformity of scattering of pathway genes along genomes are given. The results are compared with those obtained using traditional Pearson and Kolmogorov-Smirnov tests.

Keywords: randomness statistical test, Pearson and Kolmogorov-Smirnov tests

1 Introduction

Testing the uniformity of observations is a common problem in a great variety of fields such as history (distribution of specific events along the time axis), biology (distribution of special bio-markers along a chromosome), geology (time series of earthquakes), queuing theory (arrival events for the needs of customer service), evolution (time series of significant changes). The results of these tests can give a new view of the field being a good argument against “full randomness” (uniformity) while elucidating the non-uniform structure of events. In particular, overall genes are known to be evenly (uniformly) distributed along the two strands of chromosomes, while genes controlling a specific function tend to cluster on the chromosomes to facilitate co-transcription or to provide stoichiometry [2,3,4,5].

The goal of this study Our goal here is to testing the hypothesis alternative to uniformity, hypothesis (i.e. the tendency to clustering) against uniformity. For this purpose, two traditional approaches have been traditionally used, namely, the Pearson test and the Kolmogorov-Smirnov test (see, for example, [1]). Both of these give a good solution under the condition that the number of observations is large enough since because all the mathematical test results for the tests are about of the order of their asymptotical values. That is why our goal is to fill the gap between the case of large datasets, where the tools for mathematical analysis have been developed, and cases with comparatively low number of observations (<50), but with exact calculations of test characteristics such as significance or p-value. The idea of our approach is that, for a uniform distribution, the sum of squared distances between neighboring observations after their ordering should be small as compared to that for the alternative situation of clustering. In the section Method we give the formal definition of the test, the related theoretical results and a

description of its numerical evaluation. In the Application section we compare the proposed test with those, known tests from the literature, that could be applied to evaluating the distribution of the same pathway genes along the *E. coli* K-12 MG1655 chromosome belonging to same pathways. The obtained results suggest the advantage of our test in comparison with those of Pearson and Kolmogorov-Smirnov.

2 Method

Let x_1, x_2, \dots, x_n be independent observations of a random variable X with unknown continuous distribution function $F(x)$, all the observations belonging to the interval $(0,1)$. Our goal is checking the hypothesis

$$H_0 : F(x) = x, \quad 0 \leq x \leq 1,$$

i.e., the uniform of the distribution under consideration against the alternative of all other distributions. For the sake of notation simplicity, it is assumed that the observations x_1, x_2, \dots, x_n are ordered, i.e.,

$$0 = x_0 < x_1 < x_2 < \dots < x_n < x_{n+1} = 1 .$$

Statistics $S_n = (n+1) \sum_{i=0}^n \Delta_i^2$, $\Delta_i = x_{i+1} - x_i$ is introduced.

First of all, we will show that statistics S_n satisfies inequality $1 \leq S_n \leq n+1$.

Using Lagrange multipliers for the minimization of function S_n with the

restriction $\sum_{i=0}^n \Delta_i = 1$, one can easily see that $1 \leq S_n$, and from another

side $\sum_{i=0}^n \Delta_i^2 \leq (\sum_{i=0}^n \Delta_i)^2 = 1$. Thus the minimum is reached by observations

scattered along the interval with equal distances of $\frac{1}{n+1}$ between ordered

observations, while the maximum is a result of observations concentrated only at the points 0 and 1. It is shown, thus, that statistics S_n can, indeed, discriminate between two contrary levels of uniformity.

Statistical characteristics of the distribution of S_n :

The random vector $\{\Delta_0, \Delta_1, \dots, \Delta_n\}$ follows the Dirichlet distribution [1], i.e., the

uniform distribution over the simplex $\sum_{i=0}^n \Delta_i \leq 1 \subset \mathbf{R}_+^{n+1}$. Using standard

integration, one can easily get

$$E(\Delta_i) = \frac{1}{n+1}, \quad E(\Delta_i^2) = \frac{2}{(n+1)(n+2)},$$

$$Var(\Delta_i) = E(\Delta_i^2) - E^2(\Delta_i) = \frac{n}{(n+1)^2(n+2)}, \quad i = 0, \dots, n,$$

and, therefore,

$$E(S_n) = \frac{2(n+1)}{(n+2)} \rightarrow 2, \quad n \rightarrow \infty.$$

Straightforward integration gives

$$E(\Delta_i^4) = \frac{24}{(n+1)(n+2)(n+3)(n+4)},$$

$$E(\Delta_i^2 \Delta_j^2) = \frac{4}{(n+1)(n+2)(n+3)(n+4)} \quad \text{for } i \neq j,$$

$$Var(\Delta_i^2) = E(\Delta_i^4) - E^2(\Delta_i^2) = \frac{4n(5n+11)}{(n+1)^2(n+2)^2(n+3)(n+4)},$$

and, finally,

$$Var(S_n) = \frac{4n^3(n^2+5n+10)}{(n+1)^2(n+2)^2(n+3)(n+4)} = O(n^{-1}).$$

Thus for large values of n , the distribution of statistics S_n converges to the distribution concentrated at the point 2.

For numerical evaluation of the distribution function of $\frac{S_n}{n+1}$, we have developed an algorithm based on the following consideration. By definition,

$$F_n(S^2) = P\left(\sum_{i=0}^n \Delta_i^2 < S^2\right) = \int_0^s \left(\int_{\substack{\sum_{i=0}^{n-1} \Delta_i^2 < S^2 - \Delta_n^2 \\ \sum_{i=0}^{n-1} \Delta_i \leq 1 - \Delta_n}} d\Delta_0 \dots d\Delta_{n-1} \right) d\Delta_n$$

Introducing new variables $y_i = \Delta_i(1 - x_n)^{-1}$, $0 \leq i \leq n-1$, we obtain

$$F_n(S^2) = \int_0^s (1 - \Delta)^n \left(\int_{\substack{\sum_{i=0}^{n-1} y_i^2 < \frac{S^2 - \Delta^2}{(1 - \Delta_n)^2} \\ \sum_{i=0}^{n-1} y_i \leq 1}} d\Delta_0 \dots d\Delta_{n-1} \right) d\Delta =$$

$$= \int_0^s (1 - \Delta)^n F_{n-1} \left(\frac{S^2 - \Delta^2}{(1 - \Delta)^2} \right) d\Delta \quad (1)$$

The last formula gives the recursion for calculating the distribution function $F_n(S^2)$ through integrating the distribution function with $n-1$ observations

$F_{n-1} \left(\frac{S^2 - \Delta^2}{(1 - \Delta)^2} \right)$. The starting function $F_1(S^2)$ is the distribution function

for $X^2 + (1 - X)^2$, where X follows the uniform distribution on $(0,1)$. It can be

easily shown that $F_1(S^2) = 2\sqrt{0.5 \cdot S^2 - 0.25}$ for $0.5 < S^2 < 1$. The recursion

was implemented by means of consecutive numerical integration of (1) with 10000 nodes in the interval and the results are presented in Figure 1.

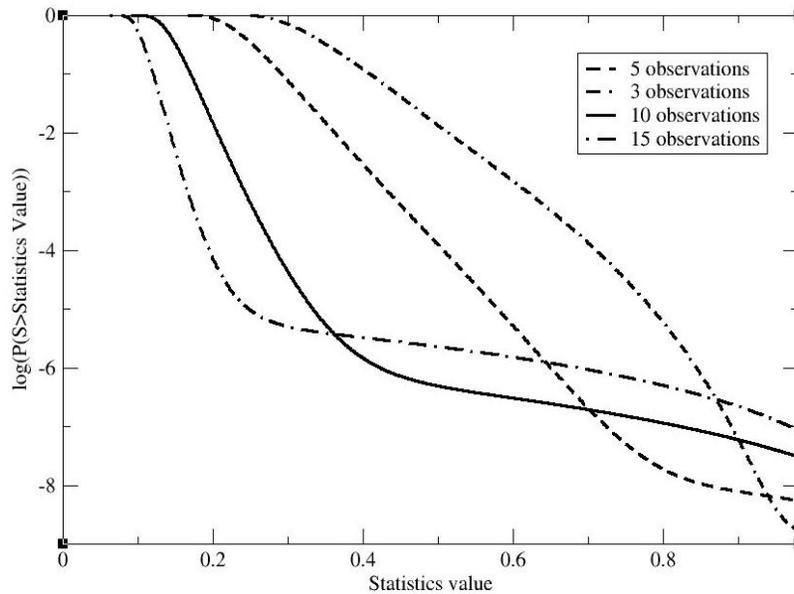


Fig. 1. Results of recursive calculations of S_n distribution, $n=3,5,10,15$.

3 Application

In this section we describe the application of the proposed technique to testing the uniformity of gene scattering in two biological pathways along bacterial chromosomes. The results are compared with those obtained by means of the Pearson χ^2 -test and Kolmogorov-Smirnov test.

To find out if the genes of certain *E. coli* pathways are uniformly distributed along the chromosome or tend to cluster, we have selected two KEGG (Kyoto

Encyclopedia of Genes and Genomes) pathways [6] to study the gene locations on the chromosomes. The position of each gene is presented by its ratio to the length of the chromosome, which allows to consider the observations as following a continuous distribution on the interval (0,1).

First we examined the small KEGG, the so-called tyrosine metabolism, pathway. Tyrosine is one of the 20 amino acids found in the cell. It has a phenol group and can serve as a precursor for the synthesis of other molecules of various types (such as hormones and pigments) in the process of metabolism. In particular, the KEGG tyrosine metabolism pathway is the one to produce tyrosine. The pathway consists of only nine *E. coli* genes, three of which belong to the leading strand, whereas the rest six belong to the lagging strand. To our knowledge, there does not exist any test that could deal with such a small number of observations.

The results of testing the uniformity hypothesis are presented in Table 1. Our calculations show that the p-values are 0.0116 and 0.3958 for the lagging strand with 6 genes and for the leading strand with 4 genes, respectively. This result enables us to unambiguously reject the hypothesis of uniformity of the pathways genes on the lagging strand and accept this hypothesis for genes on the leading strand. Thus our test can be employed in the case of very small datasets, where the traditional Pearson χ^2 and Kolmogorov-Smirnov (K-S) tests are inapplicable. Moreover, it has been demonstrated that our test allows accurate calculations of the p-value.

The second pathway with the KEGG database is that of purine metabolism, which contains a total 84 *E. coli* genes, 54 of which form 39 operons and the rest can be regarded as single-gene operons. Purines – adenine (A) and guanine (G) - are two of the four nucleotides which are used in building the DNA molecule. The purine metabolism pathway consists of a series of biochemical reactions in which A and G molecules are synthesized and degraded.

Table 1. Evaluation of p-value for the testing of uniformity of genes positions using three different tests

Number of Genes	Our Test	Pearson χ^2 test (# of intervals)					K-S test
		4	5	6	7	8	
49, leading strand	.0057	.1631	.0329	.0027	.0039	.0006	.0222
35, lagging strand	.0035	.0399	.2729	.0046	.0255	.0779	.0646

The results presented in Table 1 show fluctuation of the p-value in the Pearson test with different number of intervals, while the Kolmogorov-Smirnov test provides much higher p-values as compared to our test. Taking into account that our test is the only one that is not asymptotic, we argue that it is much more reliable than others for testing uniformity.

References

1. Walpole, R. E., Myers, R. H. and Myers, S. L., Probability and Statistics for Engineers and Scientists, 7th ed., Prentice-Hall (2002).
2. Demerec, M. and Hartman, P. E., Complex Loci in Microorganisms, *Annual Review of Microbiology*, **13**, 377-406 (1959).
3. Jacob, F. and Monod, J., Genetic regulatory mechanisms in the synthesis of proteins, *J. Mol. Biol.*, **3**, 318-356 (1961).
4. Zheng, Y., Szustakowski, J. D., Fortnow, L., Roberts, R. J. and Kasif, S., Computational identification of operons in microbial genomes, *Genome Res.*, **12**, 1221-1230 (2002).
5. Yin, Y., Zhang, H., Olman, V. and Xu, Y., Genomic arrangement of bacterial operons is constrained by biological pathways encoded in the genome, *Proc. Natl. Acad. Sci. U S A*, **107**, (2010) in press.
6. Kanehisa, M., Araki, M., Goto, S., Hattori, M., Hirakawa, M. et al., KEGG for linking genomes to life and the environment, *Nucleic Acids Res.*, **36**, D 480-484 (2008).

Stochastic Clustering and Nucleation

Florentin Paladi

Department of Theoretical Physics, State University of Moldova
A. Mateevici str. 60, Chisinau MD-2009, Moldova
Email: fpaladi@usm.md

Abstract: In this paper the present developments in the physics of complex systems, in particular the structural relaxation of supercooled liquids and glasses, are discussed by using a stochastic cluster-based model. We are able to depict the impact of the interface between the nucleus considered as a cluster of a certain number of molecules and the liquid phase for the enhancement of the overall nucleation process. In general, these mathematical models describe the interactions of agents in heterogeneous populations and they are developed within the framework of the recent discussions about the gap between agent-based computational models (ABM) and stochastic analytical models. In particular, it is shown that even a relatively simple stochastic model, which appears phenomenological if it is not agent-based, can describe precisely the outcomes from multiple agent-based simulations where there is a lack of probabilistic insight and which should be long enough to equilibrate the states of large systems.

Keywords: Stochastic modeling, Phase transitions; Nucleation; ABM

1 Introduction

In terms of clusters, growth/removal is the process in which a new cluster or a free agent is introduced/removed to/from the system. Fragmentation can be defined as a process in which a cluster breaks up into isolated agents. Coagulation is when two clusters join making a single one. Addition can be described as a special type of coagulation in which a free agent already in the system is added at random to another cluster. Attachment is the process in which new incoming agents attach themselves to an existing cluster, and this allows both the system and the clusters to grow in size. Only restricted combinations of these ingredients cause the models to differ [1–3].

One purpose of this paper is to get new insights into microscopic explanations of stochastic models which may be compared with the agent-based computational models, and to bridge the gap between agent-based models and stochastic processes. The application refers to the nucleation process, a widely spread phenomenon in both nature and technology, which may be considered as a representative of the aggregation phenomena in complex systems. The recent discovery of the generation and extinction of crystal nuclei at very low temperatures [4, 5] suggests that stochastic generation of crystal nuclei would be

considered as the result of fluctuation of complex cluster structure of the supercooled liquid. Considering that the crystal nucleation is just one extreme event in the fluctuation of a clustered structure, for example, another metastable liquid phase with a different structure from the ordinary one would also be potentially nucleated in a similar procedure. Stochastic generation of crystal nuclei may thus be considered as the result of fluctuation of cluster structure of a supercooled liquid.

The role of both heterogeneity and the interface between clusters in the enhancement of nucleation rate has still to be explained. In particular, it was observed that nuclei could almost always be formed near the surface of the cluster instead of in the interior, and one factor favoring nucleation near the surface would be the greater freedom of motion and, hence, a larger nucleation probability [6]. This is surprising because it is known that the surface layers of the nuclei tend to be disordered and melt at significantly lower temperatures than their cores.

2 The Model

Let us consider N atoms which can be in 3 different states (*cluster*, *liquid* and their *interface*), and can perform 4 possible moves: *liquid* to *interface*, *interface* to *liquid*, *interface* to *cluster*, and *cluster* to *interface*. One can identify 4 different combinations denoted with probabilities $p_1 \dots p_4$. That is, drawing randomly one particle, it will be of type i with probability p_i . Let $N=1,2,\dots,\infty$ be the total number of atoms in the system, and $\{n_1, n_2, n_3, n_4\}$ are their partition into 4 subsets. Each subset can be called cluster, and the process itself – clustering. The number of

possible partitions in this case is $P(N) = \frac{1}{3!} \prod_{i=1}^3 (N+i)$, where

$n_i = \overline{0, N}$, $i = \overline{1, 4}$ and $\sum_{i=1}^4 n_i = N$. For example, in a system of $N=1000$ atoms,

$P(N)$ equals to 167668501! Such an interaction in the ABM model always involves an active agent and a passive one: the agents have preferences over their states and they can play both roles interchangeably. Accordingly, the number of repeated computer runs due to different possible partitions would be very large.

Let's consider further that each particle interacts with the entire group both as an aggressor (in terms of the Kolmogorov theory) and as a passive agent (in terms of the ABM models) as well. Then the mean π , namely the stability index, takes here the form

$$\pi = p_1(p_1 - p_3) + (p_2 + p_3)(-p_1 + p_2 + p_3 - p_4) + p_4(p_4 - p_2)$$

or, taking into account that $\sum_{i=1}^4 p_i = 1$, one can exclude one probability, for example p_4 , from the above equation:

$$\pi = p_1(p_1 - p_3) + (1 - p_1 - 2p_2 - p_3)(1 - p_1 - p_2 - p_3) + (p_2 + p_3)(-1 + 2p_2 + 2p_3).$$

One can represent the distribution of states as a three dimensional point

$$d(l, r, c) \equiv d(p_1, p_2, p_3, p_4) = \sqrt{l^2 + r^2 + c^2},$$

where the axes are labeled l , c and r : $l = p_1 - p_3$, $c = p_2 + p_3 - p_1 - p_4$, $r = p_4 - p_2$, where $l + c + r = 0$ and $d \in [0, \sqrt{2}]$. Thus different distributions of states can lead to the same point in the sphere, i.e. different microscopic partitions can generate the same result on aggregate inside a sphere around the origin. Preliminary results for the two limit cases are obviously: if all particles would show the same behavior, then $d = \sqrt{2}$ and there is a maximum stability of states in such a completely asymmetrical system, but $\pi = 0$ for a homogeneous system, $p_1 = p_2 = p_3 = p_4 = 1/4$, and for combinations such as $p_1 = p_3$ and $p_2 = p_4$ in the case of unstable steady-states. To represent a two-dimensional graph describing the stability/instability of the system, we consider some fixed probabilities. Fig. 1 depicts such dependence

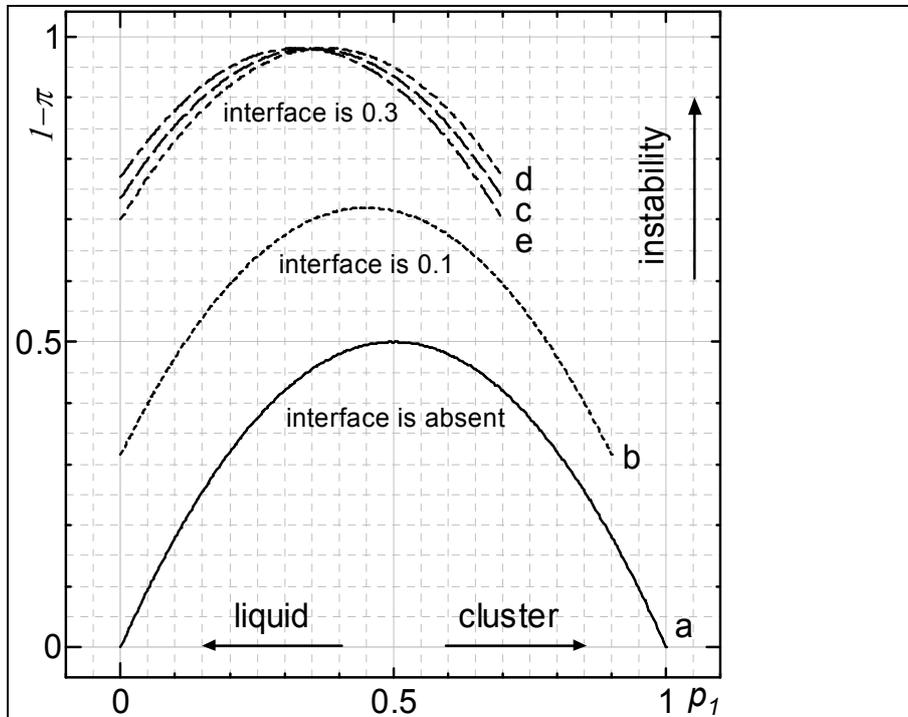


Fig. 1. Dependence of the cluster instability, $1-\pi$, on the probability p_1 for the cases: (a) particles at the liquid-cluster interface are missing; (b) share of atoms at the interface represents 1/10 of their total; corresponding share of atoms at the interface represents 3/10, i.e. $p_2=p_3=0.15$ (c), $p_2=0.1$, $p_3=0.2$ (d) and $p_2=0.2$, $p_3=0.1$ (e). A total number of 100 agents in a similar ABM model is considered.

of the cluster instability, $1-\pi$, on the probability p_1 for a system of 100 particles, where the atoms at the liquid-cluster interface are missing, (a); share of atoms at the interface p_2+p_3 is just 1/10 of their total number in the system, (b); and corresponding share of particles at the interface is equal to 3/10, i.e. $p_2=p_3=0.15$ (c), $p_2=0.1$, $p_3=0.2$ (d) and $p_2=0.2$, $p_3=0.1$ (e). Note that in the absence of an interface between liquid and cluster, as we expected, the system is in a state of maximum instability for $p_1=p_4=0.5$. While the number of particles at the interface, i.e. p_2+p_3 increases, the stability of the system decreases simultaneously, regardless of whether the particle flow at the interface is achieved at the expense of the liquid phase or particles in the cluster. However, particles at the liquid-cluster interface definitely accelerate the formation of clusters due to the displacement of the

maximum instability in the region of smaller values for p_1 . In other words, the nucleus formation is indeed a random event with a chance largely determined by the nucleation work which increases with the subnuclei size [7], and thus a decline in share of atoms, p_1 in such a cluster, namely a decline in the critical nucleus size, would be followed by the appearance of a crossover point to the supernuclei at the lower value of the energy spent on the cluster formation. Curves (d) and (e) in this figure also show that a greater flow from liquid phase ($p_2 > p_3$) or from cluster ($p_2 < p_3$) causes a minor increase in the instability of the related branches, but for $p_2 + p_3 = \text{const}$ the value of maximum π remains constant too.

3 Conclusions

We have proposed a stochastic cluster-based model for crystal nucleation. It is generally known that first-order phase transitions occur by nucleation mechanism, and both the nucleus, a cluster of molecules or atoms, and the nucleation work, a energy barrier to the phase transition, are basic thermodynamic quantities in the theory of nucleation. However, the critical nucleus formation is statistically a random event with a probability largely determined by the nucleation work. It was shown that while the number of particles at the liquid-cluster interface increases, the stability of the entire system decreases simultaneously, and the nucleus formation would be definitely enhanced due to the displacement of the bifurcation point in the region of smaller clusters. Finally, we have shown that even relatively simple stochastic models can describe precisely the results of agent-based computational models.

References

1. Rodgers, G.J., Zheng, D., *Physica A*, **308**, 375–380 (2002).
2. Rodgers, G.J., Yap, Y.J., *Eur. Phys. J. B*, **28**, 129–132 (2002).
3. Rawal, S., Rodgers, G.J., *Physica A*, **344**, 50–55 (2004).
4. Paladi, F., Oguni, M., *Phys. Rev. B*, **65**, 144202–6 (2002).
5. Paladi, F., Oguni, M., *J. Phys. Condens. Matter*, **15**, 3909–3917 (2003).
6. Chushak, Y.G., Bartell, L.S., *J. Phys. Chem. A*, **104**, 9328–9336 (2000).
7. Kashchiev, D. (2000). *Nucleation. Basic Theory with Applications*. Butterworth-Heinemann, Oxford.

A New Method for Dynamic Panel Data Models with Random Effects

by

Savas Papadopoulos

Democritus University of Thrace

Abstract

A simple-open-form estimator is introduced for the dynamic coefficient and it can be applied to levels. In a dynamic model without additional regressors, a lag of order three is included which handles the random effects. It is shown via simulations that the difference between the two OLS coefficient estimators of order-one and order-three lags estimates consistently the dynamic coefficient. In a model with other regressors, a method is suggested which estimates all the coefficients individually using restricted least squares (RLS). After estimating all the coefficients of the static regressors, the dynamic coefficient can be estimated by restricting the coefficients of the regressors to their estimated values by the suggested method (RLS) or by another method, e.g., transformed MLE (TMLE) or GMM. The RLS method can be applied when the sample size N is relatively small to the number of periods T and when the methods TMLE and GMM cannot be applied. In an application, it is shown that the RLS method provides smaller RMSE's than TMLE and GMM. Simulations compare RLS with TMLE and GMM. In general, RLS performs better than GMM. TMLE gives better results than GMM and RLS in some cases but indicates convergence problems when N and T are small.

Keywords: Restricted regression, transformed maximum likelihood, Arellano-Bond Estimator.

1. Introduction

In a regression model of panel data, we may include a lagged dependent variable as a regressor to explain dynamic economic phenomena. The appearance of a lagged dependent variable may be also used as a proxy variable for unobserved explanatory variables. In almost all economic environments, latent variables appear, e.g., quality of life, government politics, business confidence, operational risk, morale, customer satisfaction, consumer behavior, product quality, conservatism, management, marketing, etc. As a proved result in econometrics, the omission of such variables creates severe bias to their estimated coefficients, when these variables are considered as main explanatory factors. Therefore, a dynamic model would mitigate the bias of the exogenous variables included in the model. On the other hand, the exclusion of latent variables in the model would create bias in the dynamic term too. Thus, the dynamic panel data models could be applied in almost all applications with panel data, if not in all.

Probably, the most important issue in panel-data models is the individual effects. Most of the methods use first differences which theoretically wipe out those effects. Economically thinking, it makes also sense for the first differences to have individual effects to the same degree as it does for the levels. In most studies, applied researchers do not check if there are effects after taking the differences based on their theoretical elimination that could be an illusion. The application with financial ratios that we illustrate in this paper is an example with random effects on a model with first differences in the model. Differences of some order could be taken in order to deduct trend and to accomplish stationarity, and not only to remove effects. On the other hand, trend and nonstationarity could be treated by other

methods, e.g., partialing out the time variable. In this paper, the following classical dynamic model is considered and is estimated by analyzing levels and not differences,

$$(1) \quad y_{i,t} = a_i + \beta y_{i,t-1} + \mathbf{x}'_{i,t} \gamma + e_{i,t}, \quad i = 1, 2, \dots, N; \quad t = 1, 2, \dots, T.$$

The most popular methods consider Model (1) and then take the differences to wipe out the effects

$$(2) \quad \Delta y_{i,t} = \beta \Delta y_{i,t-1} + \Delta \mathbf{x}'_{i,t} \gamma + \Delta e_{i,t}.$$

Then, in model (2), we have to deal with correlation in the errors and correlation between the errors, $\Delta e_{i,t}$, and $\Delta y_{i,t-1}$. The most popular methods use GMM, e.g., Arellano Bond (1991), or maximum likelihood that takes into account the difference structure (transformed maximum likelihood, TMLE, see, Hsiao, Pesaran and Tahmiscioglu (2002)).

In this paper we introduce a new technique to estimate consistently the coefficient of the dynamic term, β , by just fitting the following restricted regression,

$$(3) \quad y_{i,t} = \tau_1 y_{i,t-1} + \tau_3 y_{i,t-3} + \mathbf{x}'_{i,t} \hat{\gamma} + e_{i,t}, \quad i = 1, 2, \dots, N; \quad t = 1, 2, \dots, T.$$

In (3) the vector parameter, γ , can be estimated by a new method introduced in Section 3 or by other methods, e.g., GMM, TMLE, etc. It turns out that, $\hat{\tau}_1 - \hat{\tau}_3$, estimates consistently the coefficient of the dynamic term, β , by restricted least squares (RLS), restricted the vector parameter, γ , to its estimate from another step of our methodology or from another method.

In the next Section a detailed literature review is presented for panel data. In Section 3, we present the general notation of our model with its assumptions and an open form of the estimator. We also explain the estimation process and we give special cases as examples for an easier understanding. Results from a simulation study are presented in Section 4, which support our method against the existing methods. An application, considered in Section 5, also supports our method since we apply out-of-sample prediction to several methods and our method gives the smaller mean-square-prediction error.

2. Literature Review for Panel Data

Extensive studies in panel data started almost fifty years ago. Zellner (1962) proved that in seemingly unrelated regression we gain more efficiency if we analyze all the equations simultaneously and not each equation individually. Balestra and Nerlove (1966) estimated a dynamic model and they noted bias for the pooled OLS and the LSDV estimator, assuming a first-order time series structure for the errors. Parks (1967) for a system of regression equations with correlated error showed asymptotic robustness for the Aitken estimator. Wallace and Hussain (1969) for a two-way-random-effect model showed asymptotic equivalence between covariance estimators and the Aitken estimator considered by Zellner (1962). Nerlove (1967) considered numerically and more carefully the bias of the dynamic modes with no effects. Thereof, Nerlove (1971) proposed a two-stage estimator for a two-way-random-effect model. In the numerical results of Nerlove (1967), it is shown that the bias increases as the error correlation increases. Later, Maddala (1971) considered MLE as an estimation technique for dynamic panel data model with error components and noticed bias.

For the two-way-random-effect model, estimates were proposed for the variances of the components by Amemiya (1971). Swamy and Arora (1972) developed a new estimator for the covariance matrix of the errors and introduced a modified Aitken estimator. Fuller and Battese (1973) suggested transformations that make the errors uncorrelated with constant

variances, but such transformations produce bias to the estimates. Fuller and Battese (1974) constructed a new estimation method for the one-way and the two-way random effects model. Their model allows constant variables over cross section or time. Avery (1977) extended the seeming unrelated regression model, considered by Zellner (1962), with error components. The estimator of the latter model was proved by Baltagi (1980) to be asymptotically inefficient and an efficient estimator was introduced by Baltagi (1980) based on Amemiya's (1971) work. Baltagi (1981a) for a non-dynamic two-way-error-component model examined the performance of several tests and estimators. Baltagi (1981b) considered simultaneous equations with error components showing that the full information estimator is more efficient than other standard alternative estimators.

The bias in dynamic panel data models with fixed effects was expressed analytically by Nickell (1981) in the simple case with no exogenous variables. Consistent instrumental variable estimators at differences are suggested by Anderson and Hsiao (1981). The analysis of differences aimed to the elimination of individual effects. Hausman and Taylor (1981) studied the instrumental-variable estimator when the individual effects are correlated with some exogenous variables. Various estimates were considered and compared under different assumptions by Anderson and Hsiao (1982). Linear and non-linear multivariate model with error components via maximum likelihood were considered by Magnus (1982). For applications of a non-linear multivariate model, see Sickles (1985), and linear multivariate model, see Sickles and Taubman (1986). Estimation for a linear system of simultaneous equations by instrumental variables was considered by Amemiya and MaCurdy (1986). Autoregressive models of high order for the endogenous and the exogenous variables were considered by Holtz-Eakin, Newey and Rosen (1988), taking the differences and using instrumental variables. The estimators by Hausman and Taylor (1981) and Amemiya and MaCurdy (1986) were compared and discussed by Breusch, Mizon and Schmidt (1989). A new estimation method for a two-way-error-component model is proposed by Wansbeek and Kapteyn (1989) for unbalanced data. For unbalanced data, the ANOVA, MLE, and MIVQUE estimators are compared by Baltagi and Chang (1994) and Baltagi, Song and Jung (2002).

Generalized method of moments (GMM) was used to estimate dynamic panel data by Arellano and Bond (1991). The elimination of the individual effects is done by taking the differences or orthogonal deviations. The efficiency for dynamic panel data under GMM can be improved by adding linear and nonlinear moment conditions, see Ahn and Schmidt (1995) and their method is supported by Wansbeek and Bekker (1996). Arellano and Bover (1995) suggested a method, with predetermined variables, for efficient IV estimators for a model with random effects. Later, Blundell and Bond (1998), suggested conditions and restrictions that improve the performance of the first-difference GMM estimator, and their work was further investigated by Hahn (1999). An alternative estimator, derived on a two-stage-least-square process, was also suggested by Keane and Runkle (1992) when the instruments are predetermined but not strictly exogenous. For more detailed discussions on all these methodologies see Arellano (2003). Bayesian estimation method is used by Hsiao and Tahmiscioglu (1997) and Tahmiscioglu (2001) on financial constraints and investment. A new estimator for nonstationary panel data is proposed by Phillips and Moon (1999). Because the LSDV estimator provides more efficiency than the GMM estimator, Kiviet (1995) proposes bias correction for LSDV and this problem was further examined by MacKinnon and Smith (1998), Hahn and Kuersteiner (2002), and Bun and Carree (2005). Also, Pesaran and Smith (1995) proved, for dynamic models with different coefficients over groups, that while the pooled and aggregated estimators are inconsistent, the cross-section estimator is consistent. Assuming T and N going to infinity, Hahn and Kuersteiner (2002) for fixed effects, and Alvarez and

Arellano (2003) for random effects, considered asymptotically bias corrected OLS, and GMM and LIML, respectively. In the latter paper, it is shown that, for fixed T and N going to infinity, GMM and LIML are consistent and asymptotically equivalent. Hsiao, Pesaran and Tahmiscioglu (2002) showed numerically that an MLE estimator, based on a transformed likelihood for dynamic panel data, has less bias than GMM and IV estimators. Their estimation method is used in our simulation and their results are also verified by our numerical results. For panel-data models with multifactor error structure, Pesaran (2006) considered common correlated effects estimators that give satisfactory results for small samples.

The last years the panel data models have been considered with spatial correlation in the errors, e.g., Baltagi (2006) and Kapoor, Kelejian, and Prucha (2007). On the choice of estimating discrete-dynamic-panel-data models, see Carro (2007). A comparison of standard errors in panel data was done by Petersen (2009) in financial data. And, many other studies have been conducted on panel data the last fifty years. We just pointed out some of them that had some specific impact on the area of panel data and especially on dynamic models. Definitely, there are still many other remarkable papers not mentioned in this paper, but it is not feasible to mention all of them. For further reference on panel data see Baltagi (2008), Hsiao (2003), Nerlove (2002) and Wooldridge (2002).

This paper aims to contribute to the area of dynamic panel data a novel estimation method, based on restricted regression, applicable to cases in which the existing methods cannot be applied, when T and N are both small.

3. Model and Estimation Procedure

We assume the following panel-data model with K dynamic equations with random effects.

$$(4) \quad x_{i,t}^{(k)} = \alpha_i^{(k)} + \beta_1^{(k)} x_{i,t-1}^{(k)} + \sum_{m=1}^{k-1} \gamma^{(k,m)} x_{i,t}^{(m)} + e_{i,t}^{(k)}, \quad k = 1, 2, \dots, K$$

The k -th random variable is regressed on its first lag and on all or some of the variables $x_{i,t}^{(m)}$, $m = 1, 2, \dots, K-1$. The model recognizes only one-way direction effects, as all the standard regression models and not as a system of simultaneous equations do, which include two-way direction effects. That is, for $k > m$ the m -th variable $x_{i,t}^{(m)}$ may affect the k -th variable $x_{i,t}^{(k)}$ but not vice versa. The individual effects $\alpha_i^{(k)}$ are assumed to be random. The errors $e_{i,t}^{(k)}$ are assumed to be independent over i , t and k .

We propose the following novel estimator for Model (4). Let us define

$$(5) \quad \mathbf{x}_0^{(k)} = \begin{pmatrix} \mathbf{x}_{1,0}^{(k)} \\ \vdots \\ \mathbf{x}_{I,0}^{(k)} \end{pmatrix}, \quad \mathbf{x}_{i,0}^{(k)} = \begin{pmatrix} x_{i,4}^{(k)} \\ \vdots \\ x_{i,T}^{(k)} \end{pmatrix}, \quad \text{for } i = 1, 2, \dots, I$$

$$(6) \quad \mathbf{X}^{(k)} = \begin{bmatrix} \mathbf{x}_{-1}^{(k)} & \mathbf{x}_{-3}^{(k)} \end{bmatrix}, \quad \mathbf{x}_{-1}^{(k)} = \begin{pmatrix} \mathbf{x}_{1,-1}^{(k)} \\ \vdots \\ \mathbf{x}_{I,-1}^{(k)} \end{pmatrix}, \quad \mathbf{x}_{-3}^{(k)} = \begin{pmatrix} \mathbf{x}_{1,-3}^{(k)} \\ \vdots \\ \mathbf{x}_{I,-3}^{(k)} \end{pmatrix},$$

$$\mathbf{x}_{i,-1}^{(k)} = \begin{pmatrix} \mathbf{x}_{i,3}^{(k)} \\ \vdots \\ \mathbf{x}_{i,T-1}^{(k)} \end{pmatrix}, \quad \mathbf{x}_{i,-3}^{(k)} = \begin{pmatrix} \mathbf{x}_{i,1}^{(k)} \\ \vdots \\ \mathbf{x}_{i,T-3}^{(k)} \end{pmatrix} \text{ for } k \geq 1$$

$$(7) \quad \hat{\theta}^{(k)} = \left(\mathbf{x}^{(k)'} \mathbf{x}^{(k)} \right)^{-1} \mathbf{x}^{(k)'} \mathbf{x}_0^{(k)}, \quad \text{for } k \geq 1$$

$$(8) \quad \hat{\mathbf{e}}^{(k)} = \mathbf{x}_0^{(k)} - \mathbf{x}^{(k)} \hat{\theta}^{(k)}, \quad \text{for } k \geq 1$$

$$(9) \quad \mathbf{y}^{(k,m)} = \begin{cases} \hat{\mathbf{e}}^{(1)}, & k=1, m=0 \\ \hat{\mathbf{q}}^{(m)}, & k \geq 2, m=k-1, \\ \left[\hat{\mathbf{q}}^{(m)} \quad \mathbf{x}_0^{(m+1)} \quad \dots \quad \mathbf{x}_0^{(k-1)} \right], & k \geq 3, m=k-2, k-3, \dots, 2, 1 \end{cases},$$

$$(10) \quad \hat{\lambda}_0^{(k,m)} = \left(\mathbf{y}^{(k,m)'} \mathbf{y}^{(k,m)} \right)^{-1} \mathbf{y}^{(k,m)'} \hat{\mathbf{e}}^{(k)}, \quad \text{for } k \geq 2 \text{ and } m = k-1, k-2, \dots, 2, 1$$

$$(11) \quad \mathbf{R}_1^{(k,m)} = \begin{bmatrix} \mathbf{0}_{(k-1-m) \times 1} & \mathbf{I}_{k-1-m} \end{bmatrix}, \quad \text{for } k \geq 2, m = k-2, k-3, \dots, 2, 1$$

(12)

$$\hat{\lambda}^{(k,m)} = \begin{cases} \hat{\lambda}_0^{(k,k-1)}, & \text{for } k \geq 2, m = k-1 \\ \hat{\lambda}_0^{(k,m)} + \left(\mathbf{y}^{(k,m)'} \mathbf{y}^{(k,m)} \right)^{-1} \mathbf{R}_1^{(k,m)'} \left[\mathbf{R}_1^{(k,m)} \left(\mathbf{y}^{(k,m)'} \mathbf{y}^{(k,m)} \right)^{-1} \mathbf{R}_1^{(k,m)'} \right]^{-1} \left(\hat{\lambda}^{(k,m+1)} - \mathbf{R}_1^{(k,m)} \hat{\lambda}_0^{(k,m)} \right), \\ \text{for } k \geq 2, m = k-2, k-3, \dots, 2, 1 \end{cases}$$

$$(13) \quad \mathbf{z}^{(k)} = \begin{cases} \begin{pmatrix} \mathbf{x}_{-1}^{(k)} & \mathbf{x}_{-3}^{(k)} \end{pmatrix}, & \text{for } k=1 \\ \begin{pmatrix} \mathbf{x}_{-1}^{(k)} & \mathbf{x}_{-3}^{(k)} & \mathbf{x}_0^{(1)} & \dots & \mathbf{x}_0^{(k-1)} \end{pmatrix}, & \text{for } k \geq 2 \end{cases},$$

$$(14) \quad \hat{\tau}_0^{(k)} = \left(\mathbf{z}^{(k)'} \mathbf{z}^{(k)} \right)^{-1} \mathbf{z}^{(k)'} \mathbf{x}_0^{(k)}, \quad \text{for } k \geq 1$$

$$(15) \quad \mathbf{R}_2^{(k)} = \begin{bmatrix} \mathbf{0}_{(k-1) \times 1} & \mathbf{0}_{(k-1) \times 1} & \mathbf{I}_{k-1} \end{bmatrix}, \quad \text{for } k \geq 2,$$

$$(16) \quad \hat{\tau}^{(k)} = \begin{cases} \hat{\tau}_0^{(k)}, & \text{for } k=1 \\ \hat{\tau}_0^{(k)} + \left(\mathbf{z}^{(k)'} \mathbf{z}^{(k)} \right)^{-1} \mathbf{R}_2^{(k)'} \left[\mathbf{R}_2^{(k)} \left(\mathbf{z}^{(k)'} \mathbf{z}^{(k)} \right)^{-1} \mathbf{R}_2^{(k)'} \right]^{-1} \left(\hat{\lambda}^{(k,1)} - \mathbf{R}_2^{(k)} \hat{\tau}_0^{(k)} \right), & \text{for } k \geq 2 \end{cases}$$

$$(17) \quad \hat{\mathbf{q}}^{(k)} = \mathbf{x}_0^{(k)} - \mathbf{z}^{(k)} \hat{\tau}^{(k)}, \quad \text{for } k \geq 2$$

Actually, the unknown parameters of Model (4) are estimated by the following $2K-1+K(K-1)/2$ regressions:

$$(18) \quad x_{i,t}^{(k)} = \theta_1^{(k)} x_{i,t-1}^{(k)} + \theta_3^{(k)} x_{i,t-3}^{(k)} + v_{i,t}^{(k)}, \quad k = 2, \dots, K$$

$$(19) \quad \hat{v}_{i,t}^{(k)} = \lambda^{(k,m)} \hat{q}_{i,t}^{(m)} + \sum_{j=m+1}^{k-1} \hat{\lambda}^{(k,j)} x_{i,t}^{(j)} + w_{i,t}^{(k,m)}, \quad k = 2, \dots, K, m = k-1, k-2, \dots, 2, 1$$

$$(20) \quad x_{i,t}^{(k)} = \tau_1^{(k)} x_{i,t-1}^{(k)} + \tau_3^{(k)} x_{i,t-3}^{(k)} + \sum_{m=1}^{k-1} \hat{\lambda}^{(k,m)} x_{i,t}^{(m)} + q_{i,t}^{(k)}, \quad k = 1, 2, 3, \dots, K$$

By the above regressions (18), (19), and (20), we actually estimate the following parameters in (4):

- i) by (19) for $k = 1, 2, \dots, K$ and $m = k-1, k-2, \dots, 2, 1$, it is estimated $\hat{\gamma}^{(k,m)}$ by $\hat{\lambda}^{(k,m)}$ and
- ii) by (20) for $k=1, 3, \dots, K$ it is estimated $\hat{\beta}_1^{(k)}$ by $\hat{\tau}_1^{(k)} - \hat{\tau}_3^{(k)}$.

The estimators presented in (5-17) are computed by regressions (18-20). Restricted regressions are executed in (19) and (20), while unrestricted multiple regressions are run in (18). The restrictions are imposed to the estimated parameters, with hat, being estimated by previous steps. In (18) we regress only the dynamic part and we also include the extra lag of order three. The additional lag of order three encounters the appearance of random effects in (4). Based on our simulations, when we use the lag of order three we have small bias while when we use the lag of order two the bias is not negligible.

A brief description about the estimation system (18-20) follows. The method consists of three stages executed for each of the k variables included in the model. In the first stage, we compute the residuals after regressing each variable on its dynamic term and on its lag of order three. In the second stage, we estimate the coefficient of the k -th variable on the m -th variable by regressing the residuals of the first stage on the residuals of the third stage from the previous iteration, $k-1$. The regressors also include all the variables with indices between $m+1$ and $k-1$ by restricting their coefficients equal to their estimated values from previous steps. The third stage estimates the coefficient of the dynamic term. The regressors include the regressors from the first stage plus all the variables with indices less than $k-1$, by restricting their coefficients equal to their estimated values from the second stage.

Equation (19), for each k , is regressed $k-1$ times estimating each time one of the coefficients of (4), $\gamma_m^{(k)}$, $m=k-1, k-2, \dots, 2, 1$, starting from $k-1$, and descending by 1 up to 1. For each k and m we regress the residuals from (18), $\hat{e}_{i,t}^{(k)}$, on the residuals from (20), $\hat{v}_{i,t}^{(m)}$, and on $x_{i,t}^{(m+1)}, \dots, x_{i,t}^{(k-1)}$, by restricting the coefficients of the variables $x_{i,t}^{(m+1)}, \dots, x_{i,t}^{(k-1)}$ equal to the estimates obtained by previous executions for smaller k . Also, the residuals $\hat{v}_{i,t}^{(m)}$ have been estimated by previous iterations of Equation (20) for smaller k . It turns out that the estimator $\hat{\lambda}^{(k,m)}$ estimates consistently the coefficient $\gamma^{(k,m)}$ of $x_{i,t}^{(m)}$ in (4) based on our simulation studies. The residuals $v_{i,t}^{(k)}$ and $q_{i,t}^{(m)}$ are free of dynamics and random effects and $q_{i,t}^{(m)}$ are also free of the variables $x_{i,t}^{(1)}, \dots, x_{i,t}^{(k-1)}$. The residuals $v_{i,t}^{(k)}$ and $q_{i,t}^{(m)}$ have the same units as the variables $x_{i,t}^{(k)}$ and $x_{i,t}^{(m)}$, respectively. Thus, in (19) by regressing $\hat{v}_{i,t}^{(k)}$ on $\hat{q}_{i,t}^{(m)}, x_{i,t}^{(m+1)}, \dots,$ and $x_{i,t}^{(k-1)}$, and by restricting the coefficients of $x_{i,t}^{(m+1)}, \dots, x_{i,t}^{(k-1)}$ on consistent estimates obtained by previous regressions, we estimate the direct effect of $\hat{q}_{i,t}^{(m)}$ on $\hat{v}_{i,t}^{(k)}$ which is equivalent to the direct effect of $x_{i,t}^{(m)}$ on $x_{i,t}^{(k)}$, according to structures of equations (4) and (19). If we do not restrict the coefficients of $x_{i,t}^{(m+1)}, \dots, x_{i,t}^{(k-1)}$ on consistent estimates then in (19) we will not estimate the direct effect of $\hat{q}_{i,t}^{(m)}$ on $\hat{v}_{i,t}^{(k)}$ since the variables $x_{i,t}^{(m+1)}, \dots, x_{i,t}^{(k-1)}$ also include the variable $x_{i,t}^{(m)}$ as it follows from (4). Note that the variables $x_{i,t}^{(m+1)}, \dots, x_{i,t}^{(k-1)}$ include indirect effects of $x_{i,t}^{(m)}$ on $x_{i,t}^{(k)}$. By this method, we can estimate the indirect effects in addition to the direct effects and therefore the total effects among the variables. Similarly, in Equation (20) we execute restricted regression with consistent estimates of $\gamma^{(k,m)}$, $m=1, 2, \dots, k-1$ from (19). Thus we restrict the coefficients of the variables $x_{i,t}^{(m)}$, $m=1, 2, \dots, k-1$ on $\hat{\lambda}^{(k,m)}$ and so we obtain the desirable direct dynamic effect. Also, Model (4) can be used for predicting all the variables at the period t using only the variables from the previous period, $t-1$.

To explain better the model we illustrate the model and the estimation process for $K=4$.
Model (4) in this case is:

$$(21) \quad x_{i,t}^{(4)} = \alpha_i^{(4)} + \beta_1^{(4)} x_{i,t-1}^{(4)} + \gamma^{(4,1)} x_{i,t}^{(1)} + \gamma^{(4,2)} x_{i,t}^{(2)} + \gamma^{(4,3)} x_{i,t}^{(3)} + e_{i,t}^{(4)},$$

$$(22) \quad x_{i,t}^{(3)} = \alpha_i^{(3)} + \beta_1^{(3)} x_{i,t-1}^{(3)} + \gamma^{(3,1)} x_{i,t}^{(1)} + \gamma^{(3,2)} x_{i,t}^{(2)} + e_{i,t}^{(3)},$$

$$(23) \quad x_{i,t}^{(2)} = \alpha_i^{(2)} + \beta_1^{(2)} x_{i,t-1}^{(2)} + \gamma^{(2,1)} x_{i,t}^{(1)} + e_{i,t}^{(2)},$$

$$(24) \quad x_{i,t}^{(1)} = \alpha_i^{(1)} + \beta_1^{(1)} x_{i,t-1}^{(1)} + e_{i,t}^{(1)},$$

For the above system (21-24) the estimation method (18-20) is written:

for $k=1$

$$(25) \quad x_{i,t}^{(1)} = \tau_1^{(1)} x_{i,t-1}^{(1)} + \tau_3^{(1)} x_{i,t-3}^{(1)} + q_{i,t}^{(1)}$$

for $k=2$

$$(26) \quad x_{i,t}^{(2)} = \theta_1^{(2)} x_{i,t-1}^{(2)} + \theta_3^{(2)} x_{i,t-3}^{(2)} + v_{i,t}^{(2)}$$

$$(27) \quad \hat{v}_{i,t}^{(2)} = \lambda^{(2,1)} \hat{q}_{i,t}^{(1)} + w_{i,t}^{(2,1)}$$

$$(28) \quad x_{i,t}^{(2)} = \tau_1^{(2)} x_{i,t-1}^{(2)} + \tau_3^{(2)} x_{i,t-3}^{(2)} + \hat{\lambda}^{(2,1)} x_{i,t}^{(1)} + q_{i,t}^{(2)}$$

for $k=3$

$$(29) \quad x_{i,t}^{(3)} = \theta_1^{(3)} x_{i,t-1}^{(3)} + \theta_3^{(3)} x_{i,t-3}^{(3)} + v_{i,t}^{(3)}$$

$$(30) \quad \hat{v}_{i,t}^{(3)} = \lambda^{(3,2)} \hat{q}_{i,t}^{(2)} + w_{i,t}^{(3,2)}$$

$$(31) \quad \hat{v}_{i,t}^{(3)} = \lambda^{(3,1)} \hat{q}_{i,t}^{(1)} + \hat{\lambda}^{(3,2)} x_{i,t}^{(2)} + w_{i,t}^{(3,1)}$$

$$(32) \quad x_{i,t}^{(3)} = \tau_1^{(3)} x_{i,t-1}^{(3)} + \tau_3^{(3)} x_{i,t-3}^{(3)} + \hat{\lambda}^{(3,1)} x_{i,t}^{(1)} + \hat{\lambda}^{(3,2)} x_{i,t}^{(2)} + q_{i,t}^{(3)}$$

and for $k=4$

$$(33) \quad x_{i,t}^{(4)} = \theta_1^{(4)} x_{i,t-1}^{(4)} + \theta_3^{(4)} x_{i,t-3}^{(4)} + v_{i,t}^{(4)}$$

$$(34) \quad \hat{v}_{i,t}^{(4)} = \lambda^{(4,3)} \hat{q}_{i,t}^{(3)} + w_{i,t}^{(4,3)}$$

$$(35) \quad \hat{v}_{i,t}^{(4)} = \lambda^{(4,2)} \hat{q}_{i,t}^{(2)} + \hat{\lambda}^{(4,3)} x_{i,t}^{(3)} + w_{i,t}^{(4,2)}$$

$$(36) \quad \hat{v}_{i,t}^{(4)} = \lambda^{(4,1)} \hat{q}_{i,t}^{(1)} + \hat{\lambda}^{(4,2)} x_{i,t}^{(2)} + \hat{\lambda}^{(4,3)} x_{i,t}^{(3)} + w_{i,t}^{(4,1)}$$

$$(37) \quad x_{i,t}^{(4)} = \tau_1^{(4)} x_{i,t-1}^{(4)} + \tau_3^{(4)} x_{i,t-3}^{(4)} + \hat{\lambda}^{(4,1)} x_{i,t}^{(1)} + \hat{\lambda}^{(4,2)} x_{i,t}^{(2)} + \hat{\lambda}^{(4,3)} x_{i,t}^{(3)} + q_{i,t}^{(4)}$$

.....

In Table 1 it is shown how the parameters of the true model in (21-24) are estimated by the fitted regressions (25-37).

Table 1. Estimators for the parameters of model equations (21-24) by regressions (25-37).

	Estimated parameter	Estimator
K=1	$\beta_1^{(1)}$	$\hat{\tau}_1^{(1)} - \hat{\tau}_3^{(1)}$
K=2	$\gamma^{(2,1)}$	$\hat{\lambda}^{(2,1)}$
	$\beta_1^{(2)}$	$\hat{\tau}_1^{(2)} - \hat{\tau}_3^{(2)}$
K=3	$\gamma^{(3,2)}$	$\hat{\lambda}^{(3,2)}$
	$\gamma^{(3,1)}$	$\hat{\lambda}^{(3,1)}$
	$\beta_1^{(3)}$	$\hat{\tau}_1^{(3)} - \hat{\tau}_3^{(3)}$
K=4	$\gamma^{(4,3)}$	$\hat{\lambda}^{(4,3)}$
	$\gamma^{(4,2)}$	$\hat{\lambda}^{(4,2)}$
	$\gamma^{(4,1)}$	$\hat{\lambda}^{(4,1)}$
	$\beta_1^{(4)}$	$\hat{\tau}_1^{(4)} - \hat{\tau}_3^{(4)}$

4. Simulation Results for RLS, TMLE, and GMM

Data were generated by Model (4) for $K=4$. The model equations are given analytically in Section 3 in Equations (21-24). In Table 2 results from simulation studies are reported by the method of Restricted Least Squares (RLS) described in Section 3, and by two of the existing methods, the Transformed Maximum Likelihood Estimator, (TMLE, see, Hsiao, Pesaran, and Tahmiscioglu (2002)), and by the Generalized Method of Moments, GMM (see, Arellano and Bond (1991)). Other methods such as bias correction methods will be considered in future studies but such studies use the GMM estimator to estimate the bias correction. Absolute bias and Root Mean Square Error, RMSE, for the dynamic coefficient, $\beta_1^{(4)} = \beta$, and the regression coefficients, $\gamma^{(4,1)} = \gamma^{(4,2)} = \gamma^{(4,3)} = \gamma$, from Equation (21) are reported. Note that the regression coefficients are not restricted to be the same but they are just given the same true values. For simplicity, for the three regression coefficients we report the average absolute bias and the average RMSE. To the sample sizes (T, N) we give the following values (5, 15), (5, 30), (5, 50), (10, 35), (10,40), (10, 100), and (100, 10), and for such a case we run the same model for true values of (β, γ) as (0.1, 0.3), (0.25, 0.25), and (0.4, 0.2). We also used the following true values for $\text{Var}\{e_{i,t}^{(k)}\} = 0.5$ and $\text{Var}\{\alpha_i^{(k)}\} = 0.2$. The same true values were also used to Equations (22-24). The model was replicated 1,000 times for each case.

In Table 2, in cases for (T, N) equal to (5, 15), (10, 35), and (100, 10), in which N is small relative to N , the methods TMLE and GMM cannot be applied since they invert a matrix that is not invertible in these cases. Note that “N/A” stands for “Not Applicable”. The method RLS not only applies, but also gives satisfactory results. On the best of our knowledge, we do not know any other method applicable for dynamic panel data models with random effects for such small N and T . Obviously this a main advantage of the suggested method RLS. In general

the bias of the dynamic coefficient under the GMM method is quite large and the RLS method performs much better than the GMM method in terms of the dynamic coefficient. The TMLE works better than the RLS method for $T=10$ but both methods give satisfactory results. For $T=5$ and $N=30$ the TMLE appears convergence problems in some iterations and the problems are more severe when $N=20$ and $N=25$. Note that the estimator of the RLS method has an open form given in (12) and (16). The TMLE and the GMM methods perform better than the RLS method for the regression coefficients, γ , when they can be applied. The RLS method can be combined with the TMLE or the GMM method. We can first apply the TMLE or the GMM method and then regress Equation (20) by restricted the coefficients $\lambda^{(k,m)}$ to the corresponding estimated coefficients by TMLE or GMM. We apply these combined techniques to financial data for out-of-sample predictions in the next section and they perform very well.

Table 2. Absolute Bias and Root MSE are given for parameters of Model (21) under three estimation methods, RLS, TMLE, and GMM for different sample sizes, T , and N and different sizes for the dynamic coefficient. Note that $\beta = \beta_1^{(4)}$ and that for the three γ parameters in (21) we report the average |Bias| and RMSE. (N/A=not applicable).

Case			Bias			RMSE		
T	N		RLS	TMLE	GMM	RLS	TMLE	GMM
5	15	$\beta=0.1$.0149	N/A	N/A	.2692	N/A	N/A
		$\gamma=0.3$.0177	N/A	N/A	.2159	N/A	N/A
		$\beta=0.25$.0230	N/A	N/A	.2868	N/A	N/A
		$\gamma=0.25$.0184	N/A	N/A	.2122	N/A	N/A
		$\beta=0.40$.0351	N/A	N/A	.3104	N/A	N/A
		$\gamma=0.20$.0182	N/A	N/A	.2096	N/A	N/A
5	30	$\beta=0.1$.0189	.0510	.1948	.1900	1.502	.2235
		$\gamma=0.3$.0101	.0077	.0038	.1412	.1995	.1158
		$\beta=0.25$.0161	.0066	.1401	.1793	.1154	.1708
		$\gamma=0.25$.0081	.0047	.0048	.1427	.1169	.1149
		$\beta=0.40$.0221	.0567	.2608	.2002	1.266	.2871
		$\gamma=0.20$.0112	.0061	.0072	.1400	.1885	.1165
5	50	$\beta=0.1$.0015	.0016	.0877	.1354	.0841	.1193
		$\gamma=0.3$.0069	.0035	.0049	.1119	.0901	.0889
		$\beta=0.25$.0002	.0036	.1296	.1441	.1106	.1585
		$\gamma=0.25$.0101	.0034	.0051	.1116	.0916	.0897
		$\beta=0.40$.0023	.0196	.1853	.1529	.1939	.2114
		$\gamma=0.20$.0118	.0036	.0065	.1112	.0946	.0900
10	35	$\beta=0.1$.0008	N/A	N/A	.0760	N/A	N/A
		$\gamma=0.3$.0071	N/A	N/A	.0698	N/A	N/A
		$\beta=0.25$.0005	N/A	N/A	.0813	N/A	N/A
		$\gamma=0.25$.0106	N/A	N/A	.0692	N/A	N/A
		$\beta=0.40$.0013	N/A	N/A	.0861	N/A	N/A
		$\gamma=0.20$.0125	N/A	N/A	.0686	N/A	N/A
10	40	$\beta=0.1$.0016	.0007	.0779	.0736	.0479	.0900
		$\gamma=0.3$.0070	.0015	.0005	.0646	.0589	.0588
		$\beta=0.25$.0014	.0019	.1042	.0757	.0527	.1146
		$\gamma=0.25$.0095	.0016	.0038	.0655	.0596	.0603
		$\beta=0.40$.0019	.0029	.1315	.0814	.0578	.1409
		$\gamma=0.20$.0112	.0016	.0074	.0651	.0581	.0600

10	100	$\beta=0.1$.0023	.0005	.0358	.0468	.0302	.0474
		$\gamma=0.3$.0057	.0010	.0009	.0418	.0370	.0370
		$\beta=0.25$.0035	.0003	.0506	.0494	.0326	.0612
		$\gamma=0.25$.0087	.0009	.0022	.0417	.0365	.0367
		$\beta=0.40$.0020	.0014	.0704	.0516	.0360	.0802
		$\gamma=0.20$.0111	.0013	.0040	.0420	.0361	.0368
100	10	$\beta=0.1$.0006	N/A	N/A	.0392	N/A	N/A
		$\gamma=0.3$.0061	N/A	N/A	.0436	N/A	N/A
		$\beta=0.25$.0026	N/A	N/A	.0367	N/A	N/A
		$\gamma=0.25$.0091	N/A	N/A	.0475	N/A	N/A
		$\beta=0.40$.0045	N/A	N/A	.0416	N/A	N/A
		$\gamma=0.20$.0106	N/A	N/A	.0411	N/A	N/A

5. Application

Panel data were analyzed for a ten-year period, (1995-2004, $T=10$) for 179 quoted Greek companies for which data were available for all ten years. Credit institutions and insurance companies were excluded. Non-consolidated annual data from the balance sheets were used. The following variables and model were fitted

$$EXPOA = \ln\left(\frac{\text{Total Operating Income} - \text{Operating Profit}}{\text{Assets}} + 2.2 + 10^{-5}\right)$$

$$EQUOL = \ln\left(\frac{\text{Equity}}{\text{Liabilities}} + 10^{-5}\right),$$

$$LIQOL = \ln\left(\frac{\text{Current Assets} - \text{Inventory}}{\text{Liabilities}} + 10^{-5}\right)$$

The numerator of the ratio in *EXPOA* is a result from financial operations and it is a measure for the expenses. The numerator of the ratio in *LIQOL* includes the so-called quick current assets, and it is a measure for liquidity. The denominators of all three ratios in *EQUOL*, *LIQOL* and *EXPOA* count for the size of the companies. In the ratios of *EQUOL* and *LIQOL* we divide by liabilities and not by assets because the numerators are included in the assets and not in the liabilities. The constants added to the ratios make the logarithmic quantities positive.

We fitted a dynamic panel data model with random effects under the RLS method described in Section 3, under the methods transformed MLE (TMLE) and GMM (Arellano-Bond), and under the combined methods RLS-TMLE and RLS-GMM and the results are shown in Table 3. The models were fitted on the first nine years 1995-2003 and the last year 2004 was used for out-of-sample predictions. The RMSE's for the out-of-sample predictions are presented in the last column of Table 3. We remind that the combined methods RLS-TMLE and RLS-GMM first apply the TMLE and the GMM methods to differences and then apply the restricted regression of Equation (20) to levels to estimate the dynamic coefficient.

The methods RLS-TMLE, RLS, and RLS-GMM provide smaller RMSE's than the GMM and TMLE methods and this is a numerical indication that they provide better estimates for this

particular application. By comparing the estimated coefficients of the regressors *EQUOL* and *LIQOL* under the methods RLS, GMM and TMLE we note that the estimates of RLS are close to the middle between the GMM and the TMLE estimates. For instance, the estimates for the coefficient of the variable *EQUOL* under the methods GMM and TMLE are 0.0017 and 0.0021 and the estimate under RLS is 0.0019, right in the middle. Similarly, the estimates for the coefficient of the variable *LIQOL* under the methods GMM and TMLE are 0.0055 and 0.0043 and the estimate under RLS is 0.0049, again right in the middle. The estimates of the dynamic coefficient for the methods RLS-TMLE, RLS, RLS-GMM, GMM and TMLE are 0.3849, 0.3819, 0.3781, 0.2223, and 0.3734. The coefficients of the dynamic coefficient for RLS-TMLE, RLS, and RLS-GMM are very close to each other and closer to the corresponding TMLE estimate than to the GMM estimate.

The results of the application indicate that the RLS method itself or combined with TMLE or GMM provide more accurate out-of-sample predictions and therefore better coefficient estimates than the GMM and the TMLE methods in the considered application. In this application if we would like to fit the model to the 30 largest companies that would be feasible by RLS and not by GMM and TMLE, due to small *N* relative to *T*, as we explained earlier.

Table 3. Estimates for the years 1995-2003 and Root Mean Square Errors of predicted values for the out of sample year of 2004.

Method	Estimates	RMSE
RLS-TMLE	$EXPOA_{i,t} = 0.38 \cdot EXPOA_{i,t-1} + 0.0021 \cdot EQUOL_{i,t} + 0.0043 \cdot LIQOL_{i,t} + 0.44 + 0.058 \cdot EXPOA_{i,t-3}$	0.0083
RLS	$EXPOA_{i,t} = 0.38 \cdot EXPOA_{i,t-1} + 0.0019 \cdot EQUOL_{i,t} + 0.0049 \cdot LIQOL_{i,t} + 0.44 + 0.058 \cdot EXPOA_{i,t-3}$	0.0084
RLS-GMM	$EXPOA_{i,t} = 0.38 \cdot EXPOA_{i,t-1} + 0.0017 \cdot EQUOL_{i,t} + 0.0055 \cdot LIQOL_{i,t} + 0.44 + 0.057 \cdot EXPOA_{i,t-3}$	0.0086
GMM	$\Delta EXPOA_{i,t} = 0.22 \cdot \Delta EXPOA_{i,t-1} + 0.0017 \cdot \Delta EQUOL_{i,t} + 0.0055 \cdot \Delta LIQOL_{i,t}$	0.0088
TMLE	$\Delta EXPOA_{i,t} = 0.37 \cdot \Delta EXPOA_{i,t-1} + 0.0021 \cdot \Delta EQUOL_{i,t} + 0.0043 \cdot \Delta LIQOL_{i,t}$	0.0096

REFERENCES

- AHN, S.C. and SCHMIDT, P. (1995). Efficient estimation of models for dynamic panel data. *Journal of Econometrics*, **68**, 5-27.
- AVERY, R.B. (1977). Error components and seemingly unrelated regressions. *Econometrica*, **45**, 199-209.
- ALVAREZ, J. and ARELLANO, M. (2003). The time series and cross-section asymptotics of dynamic panel data estimators. *Econometrica*, **71**, 1121-1159.
- AMEMIYA, T. (1971). The estimation of the variances in a variance components model. *International Economic Review*, **12**, 1-13.
- AMEMIYA, T. and MACURDY, T.E. (1986). Instrumental-variable estimation of an error-components model. *Econometrica*, **54**, 869-880.
- ANDERSON, T.W. and HSIAO, C. (1981). Estimation of dynamic models with error components. *Journal of the American Statistical Association*, **76**, 598-606.
- ANDERSON, T.W. and HSIAO, C. (1982). Formulation and Estimation of dynamic models using panel data. *Journal of Econometrics*, **18**, 47-82.
- ARELLANO, M. (2003). *Panel Data Econometrics*. Oxford: Oxford University Press
- ARELLANO, M. and BOND, S. (1991). Some tests of specification for panel data: Monte Carlo evidence and an application to employment equations, *Review of Economic Studies*, **58**, 277-297.
- ARELLANO, M. and BOVER, O. (1995). Another look at the instrumental variables estimation of error components models. *Journal of Econometrics*, **68**, 29-51.
- BALESTRA, P. and NERLOVE, P. (1966). Pooling cross section and time series model: The demand for natural gas. *Econometrica*, **34**, 585-612.
- BALTAGI, B.H. (1980). On seemingly unrelated regressions with error components. *Econometrica*, **48**, 1547-1551.
- BALTAGI, B.H. (1981a). Pooling: An experimental study of alternative testing and estimation procedures in a two-way error components model. *Journal of Econometrics*, **17**, 21-49
- BALTAGI, B.H. (1981b). Simultaneous equations with error components. *Journal of Econometrics*, **17**, 189-200.
- BALTAGI, B.H. (2006). Prediction in the panel data model with spatial correlation: the case of liquor. *Spatial Economic Analysis*, **1**, 175-185.
- BALTAGI, B.H. (2008). *Econometric Analysis of Panel Data*. 4th Edition, New York: John Wiley and Sons.

- BALTAGI, B. H. and CHANG, Y. (1994), Incomplete panels: A comparative study of alternative estimators for the unbalanced one-way error component regression model," *Journal of Econometrics*, **62**(2), 67-89.
- BALTAGI, B. H., SONG, S.H. and JUNG, B.C. (2002),. A comparative study of alternative estimators for the unbalanced two-way error component regression model. *Econometrics Journal*, **5**, 480-493.
- BREUSCH, T.S., MIZON, G.E. and SCHMIDT, P. (1989). Efficient estimation using panel data. *Econometrica*, **57**, 695-700.
- BLUNDELL, R. and BOND, S. (1998). Initial conditions and moment restrictions in dynamic panel data models. *Journal of Econometrics*. **87**, 115-143.
- BUN, M.J.G. and CARREE, M.A. (2005). Bias-corrected estimation in dynamic panel data models. *Journal of Business and Economic Statistics*, **23**, 200-210.
- CARRO, J.M. (2007). Estimating dynamic panel data discrete choice models with fixed effects. *Journal of Econometrics*, **140**, 503-528.
- FULLER, W.A. and BATTESE, G.E. (1973). Transformations for estimation of linear models with nested-error structure. *Journal of the American Statistical Association*, **68**, 626-632.
- FULLER, W.A. and BATTESE, G.E. (1974). Estimation of linear models with cross-error structure. *Journal of Econometrics*, **2**, 67-78.
- HAHN, J. (1999). How informative is the initial condition in the dynamic panel model with fixed effects? *Journal of Econometrics*, **93**, 309-326.
- HAHN, J. and KUERSTEINER, G. (2002). Asymptotically unbiased inference for a dynamic panel model with fixed effects when both n and T are large. *Econometrica*, **70**, 1639-1657.
- HAUSMAN, J.A. and TAYLOR, W.E. (1981). Panel data and unobservable individual effects. *Econometrica*, **49**, 1377-1398.
- HOLTZ-EAKIN, D., NEWEY, W. and ROSEN, H.S. (1988). Estimating vector autoregressions with panel data. *Econometrica*, **56**, 1371-1395.
- HSIAO, C. (2003). *Analysis of Panel Data*. New York: Cambridge University Press.
- HSIAO, C. and TAHMISIOGLU, A.K. (1997). A panel analysis of liquidity constraints and firm investment. *Journal of the American Statistical Association*, **92**, 455-465.
- HSIAO, C., PESARAN, M.H. and TAHMISIOGLU, A.K. (2002). Maximum likelihood estimation of fixed effects dynamic panel data models covering short time periods. *Journal of Econometrics*, **109**, 107-150.

- KAPOOR, M., KELEJIAN, H.H., PRUCHA, I.R. (2007) Panel data models with spatially correlated error components. *Journal of Econometrics*, **140**, 97-130.
- KEANE, M.P. and RUNKLE, D.E. (1992). On the estimation of panel-data models with serial correlation when instruments are not strictly exogenous. *Journal of Business & Economic Statistics*. **10**, 1-29.
- KIVIET, J.F. (1995). On bias, inconsistency, and efficiency of various estimators in dynamic panel data models. *Journal of Econometrics*, **68**, 58-78.
- MADDALA, G.S. (1971). The use of variance components models in pooling cross section and time series data. *Econometrica*, **39**, 341-358.
- MAGNUS, J.R. (1982). Multivariate error components analysis of linear and non-linear regression models by maximum likelihood. *Journal of Econometrics*. **19**, 239-287.
- MACKINNON, J.G. and SMITH, A.A. (1998). Approximate bias correction in econometrics. *Journal of Econometrics*. **85**, 205-230.
- NERLOVE, M. (1967). Experimental evidence on the estimation of dynamic economic relations from a time series of cross sections. *Economic Studies Quarterly*, **18**, 42-74.
- NERLOVE, M. (1971). Further evidence on the estimation of dynamic economic relations from a time series of cross sections. *Econometrica*, **39**, 359-382.
- NERLOVE, M. (2002). *Essays in Panel Data Econometrics*. United Kingdom: Cambridge University Press.
- NICKELL, S. (1981). Biases in dynamic models with fixed effects. *Econometrica*, **49**, 1417-1426.
- PARKS, R.W. (1967). Efficient Estimation of a System of Regression Equations when Disturbances Are Both Serially and Contemporaneously Correlated, *Journal of the American Statistical Association*, **62**, 500-509.
- PHILLIPS, P..C.B. and MOON, H.R. (1999). Linear regression limit theory for nonstationary panel data. *Econometrica*, **67**, 1057-1111.
- PESARAN, M.H. (2006). Estimation and inference in large heterogeneous panels with a multifactor error structure, *Econometrica*, **74**, 967-1012.
- PESARAN M.H. and SMITH, R. (1995). Estimating long-run relationships from dynamic heterogeneous panels, *Journal of Econometrics*, **68**, 79-113.
- PETERSEN. M.A. (2009). Estimating standard errors in finance panel data sets: Comparing approaches. *Review of Financial Studies*, **22**, 435-480.
- SICKLES, R.C. (1985). A nonlinear multivariate error components analysis of technology and specific factor productivity growth with an application to U.S. airlines. *Journal of Econometrics*, **27**, 61-78.

- SICKLES, R.C. and TAUBMAN, P. (1986). An analysis of the health and retirement study of the elderly. *Econometrica*, **54**, 1339-1356.
- SWAMY, P.A.V.B. and ARORA, S.S. (1972). The exact finite sample properties of the estimators of coefficients in the error components regression models. *Econometrica*, **40**, 261-275.
- TAHMISCIOGLU, A.K. (2001). Intertemporal variation in financial constraints on investment: A time-varying parameter approach using panel data. *Journal of Business and Economic Statistics*, **19**, 153-165.
- WALLACE, T.D. and HUSSAIN, A. (1969). The use of error components models in combining cross section with time series data. *Econometrica*, **37**, 55-72.
- WANSBEEK, T., and KAPTEYN, A. (1989). Estimation of the error-components model with incomplete panels. *Journal of Econometrics*, **41**, 341-361.
- WANSBEEK, T. and BEKKER, P. (1996). On IV, GMM and ML in a dynamic panel data models. *Economics Letters*, **51**, 145-152.
- WOOLDRIDGE, J.M. (2002). *Econometric analysis of cross section and panel data*. The MIT press.
- ZELLNER, A. (1962). An efficient method of estimating seemingly unrelated regressions and tests for aggregation bias. *Journal of the American Statistical Association*, **57**, 348-368.

SAVAS PAPADOPOULOS
ASSISTANT PROFESSOR
ISI ELECTED MEMBER
DEPARTMENT OF INTERNATIONAL ECONOMIC RELATIONS & DEVELOPMENT
DEMOCRITUS UNIVERSITY OF THRACE
UNIVERSITY CAMPUS
KOMOTINI, 69100 GREECE
E-MAIL: spapado@ierd.duth.gr
Phone: 306977343023
<http://savas-papadopoulos.blogspot.com/2010/01/homepage.html>

Modeling Biological Sequences and Web Navigation with a Semi Markov Chain

Aleka A. Papadopoulou

Department of Mathematics, Aristotle University of Thessaloniki,
Thessaloniki 54124, Greece
Email: apapado@math.auth.gr

Abstract: In the present the entrance probabilities and the probability distribution of the number of transitions to a state are studied to provide some answers to questions related to state occupancies for the semi Markov model. Biological sequences and Web navigation are two cases that initially seem to be different but to a certain extent they do have similarities. Two main aspects of word occurrences in biological sequences are: (a) where do they occur and (b) how many times do they occur. In Web navigation the similar questions are (a) when a node is visited and (b) how many times a node is visited. So, the theoretical results of this study are applied to model these two cases and derive distributions of word location or node occurrence and frequency of occurrences. Rewards/costs are included in the Web navigation model and analytic forms for the means, variances and moments of total interval rewards/costs are provided.

Keywords: Semi Markov chains, Entrance probabilities, State occupancies, Biological sequences, Words, Web navigation, Rewards/Costs

1 Introduction

In semi Markov processes we are sometimes concerned with the entrance of the process into a state rather than with the presence of the process in that state. Also because the semi Markov model allows a distinction between the number of time units that have passed and the number of transitions that have occurred we have the opportunity of asking the probability distribution of the number of transitions to a state that occurred through a time interval (Howard (1971)). An overview of probabilistic and statistical properties of words, as occurrences in biological sequences is provided in Reinert et al (2000). Studies of biological sequences using semi Markov models can be found in Chryssaphinou et al (2008), Barbu & Linnios (2008). Some cases of Markov and semi Markov reward models are examined in McClean (2004,2008), Papadopoulou (2001,2004). In this paper new sequences of matrices referring to probabilities of the semi Markov model are studied and an extension of the sequences given in Papadopoulou (1998, 2007) is given for two reasons: First, to provide some answers to questions referring to state occupancies i.e. when and how many times does a state appear and second to apply these results in order to answer the equivalent questions concerning biological sequences and web navigation. In Section 2, the entrance probabilities concerned with the number of transitions are defined and a closed analytic form in relation with the basic parameters is provided. Some asymptotic results are derived. Then,

the probability distribution of the number of transitions to a specific state through a time interval is defined and a study of the basic recursive equation applying geometric transforms is provided. In Sections 3 and 4 definitions and results of Section 2 are applied in biological sequences and web navigation. Finally, in Section 4 rewards/costs are included in the model and analytic forms for the means, variances and moments of the total interval rewards/costs produced by navigation, through a time interval, are provided.

2 The Semi Markov Model

In semi Markov processes we are sometimes concerned with the entrance of the process into a state rather than with the presence of the process in that state. Also because the semi Markov model allows a distinction between the number of time units that have passed and the number of transitions that have occurred we have the opportunity of asking the probability distribution of the number of transitions that occurred in a specific time interval. So, let us now consider a semi Markov chain with finite state space $S=\{1,2,\dots,N\}$, $\mathbf{P}(s)=\{p_{ij}(s)\}_{i,j \in S}$ the transition probability matrix of the imbedded Markov chain and $\mathbf{H}(m)=\{h_{ij}(m)\}_{i,j \in S}$ the holding time mass function matrix for the semi Markov chain. Let also $\mathbf{E}(k/n,s)=\{e_{ij}(k/n,s)\}_{i,j \in S}$ be the matrix where $e_{ij}(k/n,s)$ is the probability that the process which entered state i at time s will enter state j at time $s+n$ on its k -th transition concerning the interval $(s,s+n]$. Using probabilistic argument we can derive the following equation

$$\mathbf{E}(k/n,s) = \delta(k)\delta(n)\mathbf{I} + \sum_{m=1}^n \mathbf{C}(s,m)\mathbf{E}(k-1/n-m,s+m) \quad (1)$$

where $\mathbf{C}(s,m) = \mathbf{P}(s) \diamond \mathbf{H}(m)$ is the Hadamard product of the matrices $\mathbf{P}(s)$, $\mathbf{H}(m)$ and $\delta(n) = 1$ if $n=0$ or else $\delta(n) = 0$. If we follow a similar methodology with that of Papadopoulou (1998) we can provide a closed analytic form, in relation with the basic parameters, for the matrix $\mathbf{E}(k/n,s)$ given below

$$\begin{aligned} \mathbf{E}(0/n,s) &= \delta(n)\mathbf{I}, \quad \mathbf{E}(1/n,s) = \mathbf{C}(s,n), \quad \mathbf{E}(2/n,s) = \sum_{m=1}^n \mathbf{C}(s,m)\mathbf{C}(s+m,n-m), \\ \mathbf{E}(k/n,s) &= \sum_{j=2}^k \mathbf{S}_j(k-2,s,m_{k-2})\mathbf{C}(s+j-1,n-j+1) = \mathbf{S}_{n+1}(k-1,s,m_{k-1}), \end{aligned} \quad (2)$$

for every $k \geq 3$,

$$\mathbf{S}_j(k,s,m_k) = \sum_{m_k=2}^{j-k} \sum_{m_{k-1}=1+m_k}^{j-k+1} \cdots \sum_{m_1=1+m_2}^{j-1} \prod_{r=1}^{k-1} \mathbf{C}(s+m_{k-r}-1, m_{k-r-1}-m_{k-r}),$$

for $j \geq k + 2$, while for $j < k + 2$, $\mathbf{S}_j(k, s, m_k) = 0$.

Asymptotically and if we consider that the imbedded Markov chain converges as $s \rightarrow \infty$, i.e. $\lim_{s \rightarrow \infty} \mathbf{P}(s) = \mathbf{P}$ and take geometric transforms over k and n in equation

$$(1) \text{ we finally get that } \mathbf{E}^{gg}(y/z) = (\mathbf{I} - y\mathbf{C}^g(z))^{-1}. \quad (3)$$

Another interesting question is the one that refers to the number of transitions to a state during a time interval. Thus, if we define the matrix $\mathbf{VS}(x/n, s) = \{vs_{ij}(x/n, s)\}_{i, j \in \mathcal{S}}$, where $vs_{ij}(x/n, s)$ is the probability that the number of transitions during the interval $(s, s+n]$ to state j equals x given that the process entered state i at time s , then we can derive the following results

Case 1 In that case when $i=j$ we shall not count the initial occupancy at time s in computing the number of visits to state j . If we define as ${}^>\mathbf{W}(n, s) = \{{}^>w_i(n, s)\}_{i \in \mathcal{S}}$,

$$\text{and } {}^>w_i(n, s) = \sum_{m=n+1}^{\infty} \sum_{j=1}^N p_{ij}(s) h_{ij}(m) \text{ we can derive equation (4) below}$$

$$\begin{aligned} \mathbf{VS}(x/n, s) &= \delta(x) {}^>\mathbf{W}(n, s) \mathbf{U} + \sum_{m=0}^n \mathbf{C}(s, m) [\mathbf{VS}(x-1/n-m, s+m) \diamond \mathbf{I}] \\ &\quad + \sum_{m=0}^n \mathbf{C}(s, m) [\mathbf{VS}(x/n-m, s+m) \diamond (\mathbf{U}-\mathbf{I})] \end{aligned} \quad (4)$$

where \mathbf{U} is the matrix with all elements equal to 1. Asymptotically and if we consider that the imbedded Markov chain converges as $s \rightarrow \infty$, i.e. $\lim_{s \rightarrow \infty} \mathbf{P}(s) = \mathbf{P}$ and take geometric transforms over x and n in equation (4) we have

$$\mathbf{VS}^{gg}(y/z) = {}^>\mathbf{W}^g(z) \mathbf{U} + (y-1)\mathbf{C}^g(z) [\mathbf{VS}^{gg}(y/z) \diamond \mathbf{I}] + \mathbf{C}^g(z) \mathbf{VS}^{gg}(y/z). \quad (5)$$

Equation (5) is of the form $\mathbf{A} = \mathbf{C}_1 + \mathbf{C}_2[\mathbf{A} \diamond \mathbf{I}]$. We can use the property $[\mathbf{C}[\mathbf{A} \diamond \mathbf{I}]] \diamond \mathbf{I} = [\mathbf{C} \diamond \mathbf{I}][\mathbf{A} \diamond \mathbf{I}]$ (Howard 1971) to replace the term $[\mathbf{VS}^{gg}(y/z) \diamond \mathbf{I}]$ and then solve to find $\mathbf{VS}^{gg}(y/z)$. If we apply some more properties of the core matrix (Howard 1971) we can get the solution of equation (5) as follows

$$\mathbf{VS}^{gg}(y/z) = \frac{1}{1-z} \mathbf{U} - \frac{1-y}{1-z} [\mathbf{I} - \mathbf{C}^g(z)]^{-1} \mathbf{C}^g(z) \left[y\mathbf{I} + (1-y)[[\mathbf{I} - \mathbf{C}^g(z)]^{-1} \diamond \mathbf{I}] \right]^{-1}. \quad (6)$$

Case 2. In that case when $i=j$ we shall count the initial occupancy at time s in computing the number of visits to state j . Then we have

$$\mathbf{VS}(x/n, s) = [\mathbf{U}-\mathbf{I}] \diamond [\delta(x) {}^>\mathbf{W}(n, s) \mathbf{U} + \sum_{m=0}^n \mathbf{C}(s, m) [\mathbf{VS}(x-1/n-m, s+m) \diamond \mathbf{I}] +$$

$$\begin{aligned}
& + \sum_{m=0}^n \mathbf{C}(s, m) [\mathbf{V}\mathbf{S}(x/n - m, s + m) \diamond (\mathbf{U} - \mathbf{I})] + \delta(x-1)\delta(n) \mathbf{W}(n, s) + \\
& + \mathbf{I} \diamond \left[\sum_{m=0}^n \mathbf{C}(s, m) [\mathbf{V}\mathbf{S}(x-1/n - m, s + m) \diamond (\mathbf{U} - \mathbf{I})] \right]. \quad (7)
\end{aligned}$$

Asymptotically and if we consider that the imbedded Markov chain converges as $s \rightarrow \infty$, ($\lim_{s \rightarrow \infty} \mathbf{P}(s) = \mathbf{P}$) and take geometric transforms over x and n in (7) we get

$$\begin{aligned}
\mathbf{V}\mathbf{S}^{gg}(y/z) &= [\mathbf{I} - \mathbf{C}^g(z)]^{-1} [\mathbf{W}^g(z) [\mathbf{U} - \mathbf{I}] + y] + (y-1) [\mathbf{I} - \mathbf{C}^g(z)]^{-1} \\
& [\mathbf{C}^g(z) \mathbf{V}\mathbf{S}^{gg}(y/z) \diamond \mathbf{I}] + (y-1) [\mathbf{I} - \mathbf{C}^g(z)]^{-1} [\mathbf{C}^g(z) \diamond (\mathbf{U} - \mathbf{I})] [\mathbf{V}\mathbf{S}^{gg}(y/z) \diamond \mathbf{I}] \quad (8)
\end{aligned}$$

Equation (8) is of the form $\mathbf{A} = \mathbf{C}_1 + \mathbf{C}_2 [[\mathbf{C}_3 \mathbf{A}] \diamond \mathbf{I}] + \mathbf{C}_4 [\mathbf{A} \diamond \mathbf{I}]$. We can use again the property $[\mathbf{C}[\mathbf{A} \diamond \mathbf{I}]] \diamond \mathbf{I} = [\mathbf{C} \diamond \mathbf{I}][\mathbf{A} \diamond \mathbf{I}]$ in order to replace at first the term $[[\mathbf{C}_3 \mathbf{A}] \diamond \mathbf{I}]$ by $[\mathbf{I} - [\mathbf{C}_3 \mathbf{C}_2] \diamond \mathbf{I}]^{-1} [[\mathbf{C}_3 \mathbf{C}_1] \diamond \mathbf{I}] + [\mathbf{I} - [\mathbf{C}_3 \mathbf{C}_2] \diamond \mathbf{I}]^{-1} [[\mathbf{C}_3 \mathbf{C}_4] \diamond \mathbf{I}][\mathbf{A} \diamond \mathbf{I}]$ and result to the equation $\mathbf{A} = \mathbf{C}_1 + \mathbf{C}_2 [[\mathbf{I} - [\mathbf{C}_3 \mathbf{C}_2] \diamond \mathbf{I}]^{-1} [[\mathbf{C}_3 \mathbf{C}_1] \diamond \mathbf{I}] + [\mathbf{I} - [\mathbf{C}_3 \mathbf{C}_2] \diamond \mathbf{I}]^{-1} [[\mathbf{C}_3 \mathbf{C}_4] \diamond \mathbf{I}][\mathbf{A} \diamond \mathbf{I}]] + \mathbf{C}_4 [\mathbf{A} \diamond \mathbf{I}]$. We can apply once more the same technique, to replace the remaining term $[\mathbf{A} \diamond \mathbf{I}]$. Finally, and if we use again properties of the core matrix and equation

$$\begin{aligned}
\mathbf{W}^g(z) &= [1 - z]^{-1} [\mathbf{I} \diamond [\mathbf{C}^g(z) \mathbf{U}]] \quad \text{we can provide the solution of (8) below} \\
\mathbf{A} &= \mathbf{C}_1 + \mathbf{C}_2 [\mathbf{I} - [\mathbf{C}_3 \mathbf{C}_2] \diamond \mathbf{I}]^{-1} [[\mathbf{C}_3 \mathbf{C}_1] \diamond \mathbf{I}] + [\mathbf{C}_2 [\mathbf{I} - [\mathbf{C}_3 \mathbf{C}_2] \diamond \mathbf{I}]^{-1} [[\mathbf{C}_3 \mathbf{C}_4] \diamond \mathbf{I}] + \mathbf{C}_4] [\mathbf{I} - [\mathbf{C}_2 \diamond \mathbf{I}]] \\
& [\mathbf{I} - [\mathbf{C}_3 \mathbf{C}_2] \diamond \mathbf{I}]^{-1} [[\mathbf{C}_3 \mathbf{C}_4] \diamond \mathbf{I}] + [\mathbf{C}_4 \diamond \mathbf{I}]]^{-1} [[\mathbf{C}_1 \diamond \mathbf{I}] + [\mathbf{C}_2 \diamond \mathbf{I}][\mathbf{I} - [\mathbf{C}_3 \mathbf{C}_2] \diamond \mathbf{I}]^{-1} [[\mathbf{C}_3 \mathbf{C}_1] \diamond \mathbf{I}]], \quad (9)
\end{aligned}$$

where $\mathbf{C}_1 = [1 - z]^{-1} [\mathbf{U} - [\mathbf{I} - \mathbf{C}^g(z)]^{-1} [\mathbf{I} \diamond [\mathbf{C}^g(z) \mathbf{U}]]] + y [\mathbf{I} - \mathbf{C}^g(z)]^{-1}$,

$$\mathbf{C}_2 = (y-1) [\mathbf{I} - \mathbf{C}^g(z)]^{-1}, \quad \mathbf{C}_3 = \mathbf{C}^g(z), \quad \mathbf{C}_4 = (y-1) [\mathbf{I} - \mathbf{C}^g(z)]^{-1} [\mathbf{C}^g(z) \diamond [\mathbf{U} - \mathbf{I}]].$$

Remark In the above we had to deal with two equations of the following form $\mathbf{A} = \mathbf{C}_1 + \mathbf{C}_2 [\mathbf{A} \diamond \mathbf{I}]$, $\mathbf{A} = \mathbf{C}_1 + \mathbf{C}_2 [[\mathbf{C}_3 \mathbf{A}] \diamond \mathbf{I}] + \mathbf{C}_4 [\mathbf{A} \diamond \mathbf{I}]$. So, it is interesting to provide the general type for equations of this form and the basic steps of the applied technique in order to find the solution. The form is as follows

$$\mathbf{A} = \mathbf{C}_0 + \sum_{k=1}^n \mathbf{C}_k \left[\left[\prod_{i=0}^{k-1} \mathbf{B}_i \mathbf{A} \right] \diamond \mathbf{I} \right], \quad \mathbf{B}_0 = \mathbf{I}. \quad (10)$$

Step 1: We construct the term $\left[\left[\prod_{i=0}^{n-1} \mathbf{B}_i \mathbf{A} \right] \diamond \mathbf{I} \right]$ using equation (10) and the property

$[\mathbf{C}[\mathbf{A} \diamond \mathbf{I}]] \diamond \mathbf{I} = [\mathbf{C} \diamond \mathbf{I}][\mathbf{A} \diamond \mathbf{I}]$ and then find it in relation with the rest of the terms i.e.

$\left[\left[\prod_{i=0}^{k-1} \mathbf{B}_i \mathbf{A} \right] \diamond \mathbf{I} \right]$, $k=1, 2, \dots, n-1$. We replace the result in equation (10) and reduce it

$$\text{to } \mathbf{A} = \mathbf{C}_0^1 + \sum_{k=1}^{n-1} \mathbf{C}_k^1 \left[\prod_{i=0}^{k-1} \mathbf{B}_i \mathbf{A} \right] \diamond \mathbf{I}, \quad \mathbf{B}_0 = \mathbf{I} \quad (11)$$

$$\text{where } \mathbf{C}_0^1 = \mathbf{C}_0 + \mathbf{C}_n \left[\mathbf{I} - \prod_{i=0}^{n-1} \mathbf{B}_i \mathbf{C}_n \right]^{-1} \left[\prod_{i=0}^{n-1} \mathbf{B}_i \mathbf{C}_0 \right] \diamond \mathbf{I},$$

$$\mathbf{C}_k^1 = \mathbf{C}_k + \mathbf{C}_n \left[\mathbf{I} - \prod_{i=0}^{n-1} \mathbf{B}_i \mathbf{C}_n \right]^{-1} \left[\prod_{i=0}^{n-1} \mathbf{B}_i \mathbf{C}_k \right] \diamond \mathbf{I}.$$

Step x ($x=2, \dots, n-1$) We construct the term $\left[\prod_{i=0}^{n-x} \mathbf{B}_i \mathbf{A} \right] \diamond \mathbf{I}$ following similar

procedure as in *Step 1* and result to the equation

$$\mathbf{A} = \mathbf{C}_0^x + \sum_{k=1}^{n-x} \mathbf{C}_k^x \left[\prod_{i=0}^{k-1} \mathbf{B}_i \mathbf{A} \right] \diamond \mathbf{I}, \quad \mathbf{B}_0 = \mathbf{I} \quad (12)$$

Step n The result from *Step n-1* is the equation $\mathbf{A} = \mathbf{C}_0^{n-1} + \mathbf{C}_1^{n-1} [\mathbf{A} \diamond \mathbf{I}]$. If we follow once more similar procedure as in *Step 1* we will result to the solution $\mathbf{A} = \mathbf{C}_0^n$.

3 Modeling a biological Sequence

In what follows, the above definitions and results will be applied in biological sequences. In the present, a biological sequence is either a DNA or a protein sequence i.e. a sequence of letters either in the 4-letter DNA alphabet $\{A, C, G, T\}$ or the 20-letter amino acid alphabet. To model such a sequence we will consider the semi Markov chain with discrete finite state space $S = \{w_1, w_2, \dots, w_N\}$ where w_i , $i=1, 2, \dots, N$ is a specific word i.e. a combination of letters taken from the alphabet with known length (l_i). Through out the present we will consider only finite words and non-overlapping occurrences of them. Two main aspects of word occurrences in biological sequences are: (a) where do they occur and (b) how many times do they occur. To provide some answers, we will use the previously defined semi Markov model to derive distributions of word location and frequency of occurrences. Let us define as $\mathbf{P}(s)$ the transition probability matrix with elements equal to the probabilities of transition between the words i.e. $p_{ij}(s) = P[\text{next occurrence is of word } w_j \text{ given that the previous occurrence was of the word } w_i \text{ at position } s]$. Let us now assume that the letters that appear between successive words correspond to the holding times of the semi Markov chain. Thus, if the previous word occurred is w_i at position s , the next one is w_j and the number of

letters in between is m then the holding time in state w_i is m while the position of w_j is at $s+m+l_i$ where l_i is the length of w_i . Then, if we define as $e_{ij}(k/n, s, l_i) = P$ [the word w_j will occur at $s+n+l_i$ position and k word occurrences will happen during the interval $(s, s+n+l_i]$ given that the word w_i occurred at position s], we can derive the following recursive equation

$$e_{ij}(k/n, s, l_i) = \delta_{ij} \delta(n+l_i) \delta(k) + \sum_{k=1}^N \sum_{m=0}^n p_{ik}(s) h_{ik}(m) e_{kj}(k-1/n-m-l_k, s+m+l_i, l_k) \quad (13)$$

with initial conditions $e_{ij}(0/-l_i, s, l_i) = \delta_{ij}$, $e_{ij}(k/1-l_i, s, l_i) = 0$, $e_{ij}(k/2-l_i, s, l_i) = 0$, \dots , $e_{ij}(k/-1, s, l_i) = 0$, $e_{ij}(k/0, s, l_i) = 0$, where $h_{ij}(\cdot)$ is the conditional probability function of the number of letters that appear after the occurrence of w_i , given that the next occurrence is that of the word w_j . If we take geometric transforms over k in relation (13) and replace the transform variable with 1 we obtain interesting results concerning the probabilities of the word occurrence at specific position, i.e. probabilities concerned only with position and not with both position and number of word occurrences. Similarly, we can get recursive equations for related probabilities such as probabilities of the first occurrence of a word concerned with position or number of word occurrences or both. Finally, if we define $vs_{ij}(x/n, s, l_i)$ to be the probability that the number of occurrences of w_j , during $(s, s+n+l_i]$ equals x given that the word w_i occurred at position s , and l_i is the length of w_i , we have

$$vs_{ij}(x/n, s, l_i) = \delta(x) \delta_{ij} w_i(n, s) + \sum_{m=0}^n c_{ij}(s, m) vs_{jj}(x-1/n-m-l_j, s+m+l_i, l_j) + \sum_{k \neq j, m=0}^n c_{ik}(s, m) vs_{kj}(x/n-m-l_k, s+m+l_i, l_k). \quad (14)$$

4 Modeling Web navigation

As in section 3 if we apply definitions and results of section 2 we can model Web navigation as a semi Markov chain. The state space $S = \{1, 2, \dots, N\}$ of the chain represents the nodes that a web user possibly visits at some time. The matrix $\mathbf{P}(s)$ defines transition probabilities between the nodes and the matrix $\mathbf{H}(m)$ defines the probabilities of the holding times to the nodes. Two main aspects referring to nodes in Web navigation are (a) when a node is visited and (b) how many times a node is visited. Let $\mathbf{E}(k/n, s) = \{e_{ij}(k/n, s)\}_{i, j \in S}$ be the matrix where $e_{ij}(k/n, s)$ is the probability that the user which entered node i at time s will enter node j at time $s+n$ on its k -th transition concerning the interval $(s, s+n]$. Also let $\mathbf{VS}(x/n, s) = \{vs_{ij}(x/n, s)\}_{i, j \in S}$ be the matrix where $vs_{ij}(x/n, s)$ is the probability that the number of visits to node j , during the interval $(s, s+n]$, equals x given that the user entered node i at time s . Equations (2), (3) and (4), (6), (7), (9) provide

some answers for the distribution of node occurrences and frequency of visits to a node at any time and for the steady state. Another interesting issue arises if we include rewards/costs of making a transition from one node to another or holding to the same node for some time. So, let us define as k_{ij} the reward/cost for making a transition from node i to j and c_i the reward/cost for occupying node i during a time interval of length 1. In what follows, we provide analytic forms for the means variances and moments of the total interval cost. If we define the vector $\mathbf{v}(t, n)$ of the expected costs produced by a users navigation through the interval $(t, n]$ as follows $\mathbf{v}(t, n) = \{v_i(t, n)\}_i$, $v_i(t, n) = [\text{the expected cost produced until time } n, \text{ given that navigation started from node } i \text{ at time } t]$, it can be proved that

$$\begin{aligned} \mathbf{v}(t, n) &= (\mathbf{G}(n-t)\mathbf{1}') \diamond \mathbf{c}_1(n-t) + \mathbf{b}_1(n-t) + \\ &+ \sum_{j=2}^{n-t} [\mathbf{P} \diamond \mathbf{H}(j-1) + \mathbf{E}(j-1)] [(\mathbf{G}(n-t-j+1)\mathbf{1}') \diamond \mathbf{c}_1(n-t-j+1) + \mathbf{b}_1(n-t-j+1)] \end{aligned} \quad (15)$$

where

$$\mathbf{G}(n) = \text{diag} \left\{ \sum_{j=1}^N p_{1j} \sum_{m=n+1}^{\infty} h_{ij}(m), \dots, \sum_{j=1}^N p_{kj} \sum_{m=n+1}^{\infty} h_{kj}(m) \right\}, \mathbf{c}_1(n) = [nc_1, \dots, nc_N]'$$

$$\mathbf{b}_1(n) = \sum_{m=1}^n [(\mathbf{P} \diamond \mathbf{H}(m)) \diamond \mathbf{C}\mathbf{K}_1(m)] \mathbf{1}' = \left\{ \sum_{j=1}^N p_{ij} \sum_{m=1}^n h_{ij}(m) [mc_i + k_{ij}] \right\}_{i \in S} \quad \text{and}$$

$\mathbf{E}(n) = \{e_{ij}(n)\}_{i, j \in S}$, is the matrix of the entrance probabilities (Papadopoulou (2007)). Similarly, if we define as $\mathbf{C}\mathbf{K}_x(m)$ the matrix $\mathbf{C}\mathbf{K}_x(m) = \{(mc_i + k_{ij})^x\}_{i, j \in S}$, for $x=1, 2, \dots$ and $\mathbf{C}\mathbf{K}_0(m) = \mathbf{U}$, $\mathbf{c}_x(n)$ the vector $\mathbf{c}_x(n) = \{c_i^x n^x\}_{i \in S}$, for $x=1, 2, \dots$, $\mathbf{v}^x(t, n)$ the vector $\mathbf{v}^x(t, n) = \{v_i^x(t, n)\}_{i \in S}$, $v_i^x(t, n)$ is equal to the x -th moment of cost produced until time n , given that navigation started from node i at time t , for $x=2, 3, \dots$, $\mathbf{v}_x(t, n)$ the vector $\mathbf{v}_x(t, n) = \{[v_i(t, n)]^x\}_{i \in S}$, for $x=2, 3, \dots$, while for $x=0, 1$ $\mathbf{v}_x(t, n)$ is defined to be $\mathbf{v}_0(t, n) = \mathbf{I}$ and $\mathbf{v}_1(t, n) = \mathbf{v}(t, n)$, $\mathbf{b}_x(n)$ the vector

$$\mathbf{b}_x(n) = \sum_{m=1}^n [(\mathbf{P} \diamond \mathbf{H}(m)) \diamond \mathbf{C}\mathbf{K}_x(m)] \mathbf{1}' = \left\{ \sum_{j=1}^N p_{ij} \sum_{m=1}^n h_{ij}(m) [mc_i + k_{ij}]^x \right\}_{i \in S}$$

for $x=1, 2, \dots$, it can be proved that the vector of the r -th moments of the cost through the interval $(t, n]$ is equal to

$$\mathbf{v}^r(t, n) = (\mathbf{G}(n-t)\mathbf{1}') \diamond \mathbf{c}_r(n-t) +$$

$$+ \sum_{x=0}^r \binom{r}{x} \sum_{m=1}^{n-t} [\mathbf{P} \diamond \mathbf{H}(m)] \diamond [\mathbf{C}\mathbf{K}_x(m)] \underbrace{[\mathbf{v}(t+m, n) \diamond \dots \diamond \mathbf{v}(t+m, n)]}_{r-x \text{ terms}} \quad (16)$$

Finally, if we define the vector of the variances of the cost as follows $\mathbf{var}(t, n) = \{\text{var}_i(t, n)\}_i$, $\text{var}_i(t, n) = [\text{the variance of the cost for } (t, n) \text{ given that navigation started from node } i \text{ at time } t]$ applying the previous results we can get

$$\begin{aligned} \mathbf{Var}(t, n) &= (\mathbf{G}(n-t)\mathbf{1}) \diamond \mathbf{c}_2(n-t) + \mathbf{b}_2(n-t) + 2 \sum_{m=1}^{n-t} [\mathbf{P} \diamond \mathbf{H}(m)] \diamond \mathbf{C}\mathbf{K}_1(m) \\ &+ [(\mathbf{G}(n-t+m)\mathbf{1}') \diamond \mathbf{c}_1(n-t+m) + \mathbf{b}_1(n-t+m) \\ &+ \sum_{j=2}^{n-t+m} [\mathbf{P} \diamond \mathbf{H}(j-1) + \mathbf{E}(j-1)] [(\mathbf{G}(n-t+m-j+1)\mathbf{1}') \diamond \mathbf{c}_1(n-t+m-j+1) + \\ &+ \mathbf{b}_1(n-t+m-j+1)] + \sum_{m=1}^{n-t} [\mathbf{P} \diamond \mathbf{H}(m)] [\mathbf{v}(t+m, n) \diamond \mathbf{v}(t+m, n)] - [\mathbf{v}(t, n) \diamond \mathbf{v}(t, n)] \end{aligned}$$

References

1. Barbu V.S. and N. Limnios : *Semi Markov chains and hidden semi Markov models toward applications*, Springer (2008)
2. Chryssaphinou O., M. Karaliopoulou and N. Limnios: *On discrete time semi Markov chains and applications in word occurrences* Comm. Statist. Theory methods, 37, 1306-1322 (2008)
3. Howard, R.A.: *Dynamic Probabilistic systems*, Wiley, Chichester (1971)
4. McClean S.I., A.A. Papadopoulou and G. Tsaklidis *Discrete time reward models for homogeneous semi Markov systems*, Communications in Statistics Theory and Methods Vol. 33, No.3, 623-638, (2004)
5. McClean S., A. A. Papadopoulou, G. Tsaklidis, Lalit Garg, Maria Barton, and Peter Millard, *Evaluating Strategies using Non-homogeneous Markov and semi-Markov Systems*, Proceedings of the International Workshop in Applied Probability, Université de Technologie de Compiègne France, July 7-10, (2008).
6. Papadopoulou, A. A., *Counting transitions - entrance probabilities in non homogeneous semi-Markov systems*. Applied Stoch Models Data Anal. 13:199–206, (1998).
7. Papadopoulou A.A., *Economic rewards in non homogeneous semi Markov systems*, Communications in Statistics Theory and Methods Vol. 33, No.3, 681-696, (2004).
8. Papadopoulou A. A., G. Tsaklidis, *Some reward paths in semi Markov models with stochastic selection of the transition probabilities*, Methodology and Computing in Applied Probability 9:399–411, (2007)
9. Reinert G., S. Schbath and M.S. Waterman *Probabilistic and statistical properties of words: An overview*, journal of Computational Biology, 7, 1-2, 1-46 (2000).

Direct vs. indirect sequential Monte-Carlo filters

Yohan Petetin and François Desbouvries

Telecom SudParis/CITI Department and CNRS UMR 5157,
9 rue Charles Fourier, 91011 Evry, France
yohan.petetin@telecom-sudparis.eu, francois.desbouvries@it-sudparis.eu

Abstract. We address the recursive computation of the a posteriori filtering pdf $p_{n|n}$ in a Hidden Markov Chain (HMC). Classically $p_{n|n}$ is computed via the recursion $p_{n-1|n-1} \rightarrow p_{n|n-1} \rightarrow p_{n|n}$. In this paper we explore direct, prediction-based (P -based) and smoothing-based (S -based) alternative loops for propagating $p_{n|n}$. We next address sequential Monte Carlo (SMC) implementations of these filtering paths, and compare our algorithms via simulations.

Keywords: Sequential Monte Carlo, Particle filtering, Sampling Importance Resampling.

1 Introduction

Let (\mathbf{x}, \mathbf{y}) be an HMC : $p(\mathbf{x}_{0:n}, \mathbf{y}_{0:n}) = p(\mathbf{x}_0) \prod_{i=1}^n p(\mathbf{x}_i | \mathbf{x}_{i-1}) \prod_{i=0}^n p(\mathbf{y}_i | \mathbf{x}_i)$. Let $p_{n|m}$ be a shorthand notation for $p(\mathbf{x}_n | \mathbf{y}_{0:m})$. Bayesian filtering consists in computing $p_{n|n}$, or at least some approximation of the measure $p(d\mathbf{x}_n | \mathbf{y}_{0:n})$ with pdf $p(\mathbf{x}_n | \mathbf{y}_{0:n})$. $p_{n|n}$ can be computed from $p_{n-1|n-1}$ by using the path $p_{n-1|n-1} \xrightarrow{P} p_{n|n-1} \xrightarrow{U} p_{n|n}$, in which we first predict state \mathbf{x}_n , based on the same measurements (whence superscript P), and then update the measurements set $\{\mathbf{y}_k\}_{k=0}^{n-1}$ with the new data \mathbf{y}_n (whence superscript U). This path is described by the well known equation (here \mathcal{N} stands for numerator) :

$$p(\mathbf{x}_n | \mathbf{y}_{0:n}) = \frac{p(\mathbf{y}_n | \mathbf{x}_n) \int \overbrace{p(\mathbf{x}_n | \mathbf{x}_{n-1}) p(\mathbf{x}_{n-1} | \mathbf{y}_{0:n-1})}^{p(\mathbf{x}_n | \mathbf{y}_{0:n-1})} d\mathbf{x}_{n-1}}{p(\mathbf{y}_n | \mathbf{y}_{0:n-1}) = \int \mathcal{N} d\mathbf{x}_n}. \quad (1)$$

However, computing (1) is often impossible in practice, so many approximate techniques have been developed. Among them, particle filters (PF) [5] [1] are SMC methods which propagate a discrete approximation of $p(d\mathbf{x}_n | \mathbf{y}_{0:n})$.

In this paper we do not try to further improve the PF algorithms based on (1); we rather focus on (1) itself, or indeed explore alternate paths for computing $p_{n|n}$ recursively, even if $p_{n|n}$ is obtained as a byproduct.

Let us consider only those paths in which one time index is incremented at a time. The first alternative is $p_{n-1|n-1} \rightarrow p_{n-1|n} \rightarrow p_{n|n}$. Both paths compute $p_{n|n}$ recursively and differ only by the intermediate step which is

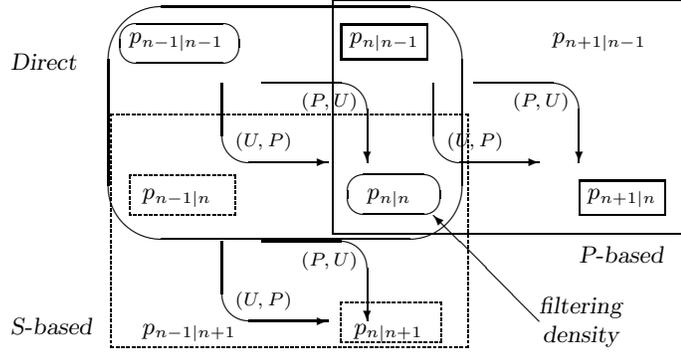


Fig. 1. The direct, P -based and S -based filtering paths.

either the one step predictive pdf $p_{n|n-1}$, or the one step smoothing pdf $p_{n-1|n}$. Now, in turn $p_{n|n-1}$ and $p_{n-1|n}$ can be propagated via the two paths obtained by moving one index and next the other. This observation yields six paths for computing $p_{n|n}$ recursively; the two paths already mentioned are "direct", i.e. $p_{n|n}$ is computed as the output of a loop with input $p_{n-1|n-1}$; two other paths are P -based, i.e. $p_{n|n}$ is computed indirectly from $p_{n|n-1}$, but the recursion itself now acts on $p_{n|n-1}$; and two paths are S -based, see Fig. 1. Out of these 6 paths only 4 are distinct, because the two paths at the boundary direct/ P -based and direct/ S -based coincide (for instance, the direct path $p_{n-1|n-1} \rightarrow p_{n|n-1} \rightarrow p_{n|n}$ coincides, up to a shift in time, with the P -based path $p_{n|n-1} \rightarrow p_{n|n} \rightarrow p_{n+1|n}$). The paper is organized as follows. We recall the four direct and indirect filtering paths in §2 and consider their SMC implementations in §3 (see [3] for details). §4 is devoted to simulations.

2 Direct, P -based and S -based paths

- The direct path $p_{n-1|n-1} \rightarrow p_{n|n-1} \rightarrow p_{n|n}$ is described by (1). Since it involves the one-step ahead prediction pdf $p_{n|n-1}$ we will call it 1- P ;
- The alternate direct path $p_{n-1|n-1} \rightarrow p_{n-1|n} \rightarrow p_{n|n}$ involves the one step backward smoothing pdf $p_{n-1|n}$ and will thus be denoted as 1- S :

$$p(\mathbf{x}_n | \mathbf{y}_{0:n}) = \int p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{y}_n) \underbrace{\left[\frac{p(\mathbf{y}_n | \mathbf{x}_{n-1}) p(\mathbf{x}_{n-1} | \mathbf{y}_{0:n-1})}{p(\mathbf{y}_n | \mathbf{y}_{0:n-1})} \right]}_{p(\mathbf{x}_{n-1} | \mathbf{y}_{0:n})} d\mathbf{x}_{n-1}; \quad (2)$$

- P -based paths compute $p_{n|n}$, but via a recursive loop involving $p_{n|n-1}$. Path $p_{n|n-1} \rightarrow p_{n|n} \rightarrow p_{n+1|n}$ coincides with 1- P (up to a shift in time). The other P -based path $p_{n|n-1} \rightarrow p_{n+1|n-1} \rightarrow p_{n+1|n}$ involves the two-

step ahead prediction pdf $p_{n+1|n-1}$ and will thus be denoted by 2- P :

$$p(\mathbf{x}_{n+1}|\mathbf{y}_{0:n}) = \frac{p(\mathbf{y}_n|\mathbf{x}_{n+1}, \mathbf{y}_{0:n-1}) \int \overbrace{p(\mathbf{x}_{n+1}|\mathbf{y}_{0:n-1})} p(\mathbf{x}_{n+1}|\mathbf{x}_n)p(\mathbf{x}_n|\mathbf{y}_{0:n-1})d\mathbf{x}_n}{p(\mathbf{y}_n|\mathbf{y}_{0:n-1}) = \int \mathcal{N}d\mathbf{x}_{n+1}}; \quad (3)$$

- S -based paths compute $p_{n|n}$, but via a recursive loop involving $p_{n-1|n}$. Path $p_{n-1|n} \rightarrow p_{n|n} \rightarrow p_{n|n+1}$ coincides with (2) (up to a shift in time). The other S -based path $p_{n-1|n} \rightarrow p_{n-1|n+1} \rightarrow p_{n|n+1}$ involves the two-step backward smoothing pdf $p_{n-1|n+1}$ and will be denoted by 2- S :

$$p(\mathbf{x}_n|\mathbf{y}_{0:n+1}) = \int p(\mathbf{x}_n|\mathbf{x}_{n-1}, \mathbf{y}_n, \mathbf{y}_{n+1}) \underbrace{\frac{p(\mathbf{y}_{n+1}|\mathbf{x}_{n-1}, \mathbf{y}_n)p(\mathbf{x}_{n-1}|\mathbf{y}_{0:n})}{p(\mathbf{y}_{n+1}|\mathbf{y}_{0:n}) = \int \mathcal{N}d\mathbf{x}_{n-1}}}_{p(\mathbf{x}_{n-1}|\mathbf{y}_{0:n+1})} d\mathbf{x}_{n-1}. \quad (4)$$

3 SMC implementations

3.1 A practical toolbox

Each path (1) to (4) is made of the succession of a propagation step P (which transforms some pdf $p(\mathbf{x}_1)$ into $p(\mathbf{x}_1) \xrightarrow{P} p(\mathbf{x}_2) = \int p(\mathbf{x}_2|\mathbf{x}_1)p(\mathbf{x}_1)d\mathbf{x}_1$), and of a Bayesian or updating step U (which transforms some density $p(\mathbf{x})$ into $p(\mathbf{x}) \xrightarrow{U} p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$), or vice versa. Let us recall how we can propagate a set of points sampled from $p(\mathbf{x}_1)$ (resp. from $p(\mathbf{x})$) into a set of points sampled (at least approximatively) from $p(\mathbf{x}_2)$ (resp. from $p(\mathbf{x}|\mathbf{y})$).

1. *Propagating*. Starting from N i.i.d. samples $\{\mathbf{x}_1^i\}_{i=1}^N \sim p(\mathbf{x}_1)$, we get N i.i.d. samples $\{\mathbf{x}_2^i\}_{i=1}^N \sim p(\mathbf{x}_2)$ by sampling, for each i , \mathbf{x}_2^i from $p(\mathbf{x}_2|\mathbf{x}_1^i)$. This is nothing but the Sampling step S of PF algorithms;
2. *Updating* [8]. Starting from $\{\mathbf{x}^i\}_{i=1}^N \sim p(\mathbf{x})$, we get N points $\{\tilde{\mathbf{x}}^i\}_{i=1}^N$ (approximately) independently sampled from $p(\mathbf{x}|\mathbf{y})$ by associating to each sample \mathbf{x}^i a weight proportional to $p(\mathbf{y}|\mathbf{x}^i)$, and then sampling $\{\tilde{\mathbf{x}}^i\}_{i=1}^N \sim \sum_{i=1}^N \frac{p(\mathbf{y}|\mathbf{x}^i)}{\sum_{i=1}^N p(\mathbf{y}|\mathbf{x}^i)} \delta_{\mathbf{x}^i}(d\mathbf{x})$. We just described nothing but the Weighting step W , followed by the Resampling step R of PF algorithms.

3.2 SMC algorithms

We now routinely derive generic SMC implementations of (1) to (4).

- 1- P . (1) gives the Bootstrap [6] : Let $p(d\mathbf{x}_{n-1}|\mathbf{y}_{0:n-1}) \approx \sum_{i=1}^N \frac{1}{N} \delta_{\mathbf{x}_{n-1}^i}(d\mathbf{x}_{n-1})$.
 S . For $1 \leq i \leq N$, sample $\tilde{\mathbf{x}}_n^i$ from $p(\mathbf{x}_n|\mathbf{x}_{n-1}^i)$;

- W. For $1 \leq i \leq N$, compute $w_n^i \propto p(\mathbf{y}_n | \tilde{\mathbf{x}}_n^i)$, $\sum_i w_n^i = 1$;
R. For $1 \leq i \leq N$, sample \mathbf{x}_n^i from $\sum_{i=1}^N w_n^i \delta_{\tilde{\mathbf{x}}_n^i}(\mathbf{d}\mathbf{x}_n)$.
- 1-S. (2) gives [2, Algorithm 8.1.1. p. 253], which is a *reordering* of the SIR algorithm with optimal importance distribution $p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{y}_n)$ (and systematic resampling) (see [4], where the successive steps are $S \rightarrow W \rightarrow R$):
Let $p(\mathbf{d}\mathbf{x}_{n-1} | \mathbf{y}_{0:n-1}) \approx \sum_{i=1}^N \frac{1}{N} \delta_{\mathbf{x}_{n-1}^i}(\mathbf{d}\mathbf{x}_{n-1})$;
W. For $1 \leq i \leq N$, compute $w_n^i \propto p(\mathbf{y}_n | \mathbf{x}_{n-1}^i)$, $\sum_{i=1}^N w_n^i = 1$;
R. For $1 \leq i \leq N$, sample $\tilde{\mathbf{x}}_{n-1}^i \sim \sum_{i=1}^N w_n^i \delta_{\mathbf{x}_{n-1}^i}(\mathbf{d}\mathbf{x}_{n-1})$;
S. For $1 \leq i \leq N$, sample \mathbf{x}_n^i from $p(\mathbf{x}_n | \tilde{\mathbf{x}}_{n-1}^i, \mathbf{y}_n)$.
- 2-P. Implementing (3) would require the knowledge of $p(\mathbf{y}_n | \mathbf{x}_{n+1}, \mathbf{y}_{0:n-1})$, but this pdf is not directly available. We thus consider the alternative path $p_{n|n-1} \rightarrow p_{n,n+1|n-1} \rightarrow p_{n,n+1|n} \rightarrow p_{n+1|n}$, given by

$$p(\mathbf{x}_{n+1} | \mathbf{y}_{0:n}) = \int \frac{\overbrace{p(\mathbf{x}_n, \mathbf{x}_{n+1} | \mathbf{y}_{0:n-1})}^{p(\mathbf{x}_n, \mathbf{x}_{n+1} | \mathbf{y}_{0:n-1})}}{p(\mathbf{y}_n | \mathbf{y}_{0:n-1})} \frac{[p(\mathbf{x}_{n+1} | \mathbf{x}_n) p(\mathbf{x}_n | \mathbf{y}_{0:n-1})]}{\int \mathcal{N} \mathbf{d}\mathbf{x}_n \mathbf{d}\mathbf{x}_{n+1}} \mathbf{d}\mathbf{x}_n. \quad (5)$$

Let us implement (5). Let $p(\mathbf{d}\mathbf{x}_n | \mathbf{y}_{0:n-1}) \approx \sum_{i=1}^N \frac{1}{N} \delta_{\mathbf{x}_n^i}(\mathbf{d}\mathbf{x}_n)$.

- S. For $1 \leq i \leq N$, sample $\tilde{\mathbf{x}}_{n+1}^i \sim p(\mathbf{x}_{n+1} | \mathbf{x}_n^i)$;
W. For $1 \leq i \leq N$, compute $w_n^i \propto p(\mathbf{y}_n | \mathbf{x}_n^i)$, $\sum_{i=1}^N w_n^i = 1$;
R. For $1 \leq i \leq N$, sample \mathbf{x}_{n+1}^i from $\sum_{i=1}^N w_n^i \delta_{\tilde{\mathbf{x}}_{n+1}^i}(\mathbf{d}\mathbf{x}_{n+1})$.

Filtering. $p(\mathbf{d}\mathbf{x}_n | \mathbf{y}_{0:n}) \approx \sum_{i=1}^N w_n^i \delta_{\mathbf{x}_n^i}(\mathbf{d}\mathbf{x}_n)$.

- 2-S. We now implement (4). Let $p(\mathbf{d}\mathbf{x}_{n-1} | \mathbf{y}_{0:n}) \approx \sum_{i=1}^N \frac{1}{N} \delta_{\mathbf{x}_{n-1}^i}(\mathbf{d}\mathbf{x}_{n-1})$.
W. For $1 \leq i \leq N$, compute $w_{n+1}^i \propto p(\mathbf{y}_{n+1} | \mathbf{x}_{n-1}^i, \mathbf{y}_n)$, $\sum_{i=1}^N w_{n+1}^i = 1$;
R. For $1 \leq i \leq N$, sample from $\sum_{i=1}^N w_{n+1}^i \delta_{\mathbf{x}_{n-1}^i}(\mathbf{d}\mathbf{x}_{n-1})$. We get N points $\tilde{\mathbf{x}}_{n-1}^i$, (approximately) distributed $\sim p(\mathbf{d}\mathbf{x}_{n-1} | \mathbf{y}_{0:n+1})$;
S. For $1 \leq i \leq N$, sample $\mathbf{x}_n^i \sim p(\mathbf{x}_n | \tilde{\mathbf{x}}_{n-1}^i, \mathbf{y}_n, \mathbf{y}_{n+1})$; then $p(\mathbf{d}\mathbf{x}_n | \mathbf{y}_{0:n+1}) \approx \sum_{i=1}^N \frac{1}{N} \delta_{\mathbf{x}_n^i}(\mathbf{d}\mathbf{x}_n)$.
Filtering. For $1 \leq i \leq N$, sample $\bar{\mathbf{x}}_{n+1}^i \sim p(\mathbf{x}_{n+1} | \mathbf{x}_n^i, \mathbf{y}_{n+1})$; then $p(\mathbf{d}\mathbf{x}_{n+1} | \mathbf{y}_{0:n+1}) \approx \sum_{i=1}^N \frac{1}{N} \delta_{\bar{\mathbf{x}}_{n+1}^i}(\mathbf{d}\mathbf{x}_{n+1})$.

4 Simulations

4.1 Simulations, linear model

We first consider the state-space model : $x_{n+1} = 0.2x_n + u_n$, $y_n = 5x_n + v_n$, in which $u_n \sim \mathcal{N}(0, Q)$ and $v_n \sim \mathcal{N}(0, R)$ are i.i.d., mutually independent and independent of $x_0 \sim \mathcal{N}(0.5, 0.5)$. Even though exact Kalman filtering (KF) is available, we compare to that benchmark solution the four SMC algorithms

2-*P*, 1-*P*, 1-*S*, 2-*S*, and the SIR algorithm with optimal importance function $p(\mathbf{x}_n|\mathbf{x}_{n-1}, \mathbf{y}_n)$ (simply denoted SIR). Let $\mathcal{J} = \frac{1}{50} \sum_{n=1}^{50} (\frac{1}{200} \sum_{j=1}^{200} (\hat{x}_{n|n}^j - x_n^j)^2)^{\frac{1}{2}}$ (200 is the number of realizations). Let $N = 50$ and $R = 2$. As we see from Table 1 *S*-based algorithms outperform *P*-based ones, and for a class of algorithms (*P*- or *S*-based) better results occur when updating precedes propagation. For $Q = 0.1$ all algorithms are similar, but the *P*-based ones degrade as Q increases, for strong variations of x_n are better tracked when y_n (for SIR or 1-*S*) or y_n and y_{n+1} (for 2-*S*) are taken into account. Next in Fig. 2 $Q = 0.1$ and \mathcal{J} evolves with N . The ordering of the algorithms is maintained, but when N increases SIR, 1-*S*, 2-*S* and KF become very close.

	2- <i>P</i>	1- <i>P</i>	SIR	1- <i>S</i>	2- <i>S</i>	KF
$Q = 0.1$	0.221996933	0.217401933	0.215515333	0.2136542	0.213415733	0.2116127
$Q = 1$	0.4304519	0.2955335	0.2745811	0.2724687	0.2723657	0.2696133
$Q = 3$	0.8114436	0.3258072	0.2827361	0.2801141	0.2773301	0.2766857
$Q = 5$	1.0114103	0.3932067	0.2856878	0.2840456	0.2834111	0.2810438
$Q = 10$	1.5200521	0.4607805	0.2870748	0.2853349	0.2848457	0.2822257

Table 1. Empirical standard deviation \mathcal{J} , linear model.

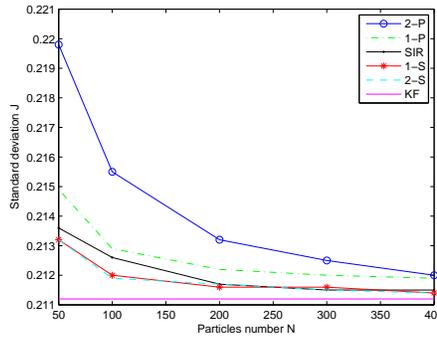


Fig. 2. Empirical standard deviation \mathcal{J} , linear model.

4.2 Simulations, Kitagawa model

Let us now consider the model

$$\begin{cases} x_{n+1} = f_n(x_n) + u_n \\ y_n = x_n^2/20 + v_n \end{cases}, \tag{6}$$

with $f_n(x_n) = 0.5x_n + 25\frac{x_n}{1+x_n^2} + 8 \cos(1.2(n+1))$, and $u_n \sim \mathcal{N}(0, Q)$ and $v_n \sim \mathcal{N}(0, R)$ are i.i.d., mutually independent and independent of $x_0 \sim \mathcal{N}(0.5, 0.5)$. In (6) $p(x_n|x_{n-1}, y_n)$ and $p(y_n|x_{n-1})$ cannot be computed exactly. So we use some approximation (linearization [4], [7], EMM [7] or UPF [9]).

	2-P	1-P	SIR(EMM)	SIR(UPF)	1-S(EMM)	1-S(UPF)
$R = 0.5$	5.0124197	3.7637961	3.4181757	3.5046666	2.7191797	2.734824
$R = 1$	4.422289	3.6087193	3.447493	3.6076351	2.8830467	2.9981435
$R = 5$	4.1079475	3.6386135	3.705746	3.6928957	3.3874559	3.3369996
$R = 10$	4.4682053	4.0435623	4.020258	4.0743315	3.7412504	3.8619037
$R = 20$	5.0894144	5.0229244	4.8433825	4.8105053	4.7580243	4.8140369

Table 2. Empirical standard deviation \mathcal{J} , Kitagawa model.

	2-P	1-P	SIR(EMM)	SIR(UPF)	1-S(EMM)	1-S(UPF)
$N = 50$	7.1328	5.5397	5.2445	5.1169	4.7151	4.8542
$N = 100$	5.9246	4.9741	4.8057	4.9631	4.6840	4.8111
$N = 150$	5.4544	4.7726	4.8290	4.7241	4.5828	4.5029
$N = 200$	5.1171	4.6771	4.5407	4.7033	4.6551	4.5327
$N = 300$	4.9559	4.5545	4.4494	4.4312	4.4449	4.4393

Table 3. Empirical standard deviation \mathcal{J} , Kitagawa model.

Let $Q = 1$, $N = 50$, and in UPF $\alpha = 1$ and $\beta = 0$. Table 2 displays \mathcal{J} for different values of R . The ordering of the algorithms is maintained; the difference between SIR and 1-S becomes significant; and for 1-S and SIR EMM provides better results than UPF. Note that in (6) f_n has strong variations, so the influence of the new data y_n is essential. This explains the difference between P -based algorithms and the SIR and 1-S ones, at least when R is small. On the other hand if the observations become very noisy ($R = 10$) the performance of 2-P is unaltered, while SIR and 1-S degrade.

Next in Table 3 we set $Q = 10$ and $R = 1$ and we see how \mathcal{J} evolves with N . Let now $\alpha = 0.94$ and $\beta = 0$. As above, 1-S outperforms the P -based and the SIR algorithms, and we observe e.g. that 1-S(EMM) with $N = 50$ particles gives about the same result as 1-P with $N = 200$ particles.

4.3 Simulations, semi-linear models

In §4.2 we compared the P -based, SIR and 1-S algorithms, but not the 2-S one, because in (6) $p(x_n|x_{n-1}, y_n, y_{n+1})$ and $p(y_{n+1}|x_{n-1}, y_n)$ are difficult to

implement. Yet 2-S can be used in some situations. Let us first consider the non linear model with linear measurements equation

$$\begin{cases} x_{n+1} = f_n(x_n) + u_n, \\ y_n = 0.5x_n + v_n, \end{cases} \tag{7}$$

in which $u_n \sim \mathcal{N}(0, Q)$ and $v_n \sim \mathcal{N}(0, R)$ are i.i.d., mutually independent and independent of $x_0 \sim \mathcal{N}(0, 1)$ (the first equation of (6) and (7) coincide).

In (7) $p(x_n|x_{n-1}, y_n)$ and $p(y_n|x_{n-1})$ can be computed easily. $p(x_n|x_{n-1}, y_n, y_{n+1})$ cannot, but the problem of computing $p(x_n|x_{n-1}, y_n, y_{n+1})$ from $(p(x_n|x_{n-1}, y_n), p(y_{n+1}|x_n))$ is the same as that of computing $p(x_n|x_{n-1}, y_n)$ from $(p(x_n|x_{n-1}), p(y_n|x_n))$ and so the approximation techniques recalled in §4.2 can be adapted to (7) (a difference however is that the exact moments of $p(x_n|x_{n-1}, y_n, y_{n+1})$ cannot be computed in (7)).

Let $R = 2$ and $N = 50$. For 2-S we use either a second-order Taylor series expansion, or UPF with $\alpha = 0.73$ and $\beta = \alpha^2 - 1$. Table 4 displays \mathcal{J} as a function of Q . As we can see, the ordering 2-P < 1-P < SIR < 1-S is maintained. 2-S outperforms 1-S if Q is small, but 1-S performs better if Q increases. The reason why is that in (7) $p(x_n|x_{n-1}, y_n)$ (used in 1-S) can be computed exactly but $p(x_n|x_{n-1}, y_n, y_{n+1})$ (used in 2-S) cannot. Since all techniques indeed approximate f_n (up to the first orders), the results strongly depend on this function. In (6) f_n has strong variations; all orders matter, so all approximations of f_n at some point becomes very poor outside of a small neighbourhood of that point, and such situations do happen if Q gets large.

	2-P	1-P	SIR	1-S	2-S(Taylor)	2-S(UPF)
$Q = 0.1$	1.649036	1.5608385	1.5513972	1.1820396	1.1017785	1.0886346
$Q = 1$	2.1070664	1.8257239	1.6739058	1.5607033	2.4151141	1.6287723
$Q = 10$	2.8134363	2.4198548	2.3303342	2.2850875	2.8866604	2.4475675
$Q = 50$	3.8612771	2.9077013	2.7306043	2.70971	2.7548658	2.7188303

Table 4. Empirical standard deviation \mathcal{J} , semi-linear model.

Finally let us consider the semi-linear model (7), but in which the evolution equation is replaced by $x_{n+1} = \arctan x_n + u_n$. Let $\alpha = 0.88$, $\beta = \alpha^2 - 1$, $N = 50$ and $R = 2$. Table 6 displays \mathcal{J} in terms of Q . By contrast with (7), function $f_n = \arctan$ is now very smooth. As a result all algorithms give satisfactory results, especially if Q is low. Also observe that the ordering of the algorithms is maintained, and in particular that 2-S outperforms 1-S, even when Q is large. The reason why is that for the arctan function limited order approximations are valid in a large domain, so the necessity of approximating $p(x_n|x_{n-1}, y_n, y_{n+1})$ is no longer a handicap of 2-S w.r.t. 1-S.

	2-P	1-P	SIR	1-S	2-S(UPF)
$Q = 1$	1.1240331	1.1202899	1.1118075	1.1041993	1.1033117
$Q = 5$	1.861067	1.8353582	1.8334354	1.8134666	1.8107098
$Q = 20$	2.611825	2.4849116	2.439016	2.4268744	2.4144096
$Q = 50$	3.1980866	2.761685	2.6592343	2.6345335	2.6329076

Table 5. Empirical standard deviation \mathcal{J} , alternate semi-linear model.

5 Conclusion

We explored direct and indirect paths (and their generic SMC implementations) for computing $p_{n|n}$ recursively. These algorithms remain PF, in the sense that their aim is to compute $p_{n|n}$, but possibly via a predictive or smoothing distribution. Our algorithms were validated by simulations. S-based algorithms outperform P-based ones, and in each class of algorithms better results are obtained (under fair conditions, i.e. when the necessary approximations are valid) when updating precedes propagation.

References

1. Arulampalam, M. S., Maskell, S., Gordon, N., and Clapp, T., "A Tutorial on Particle Filters for Online Nonlinear / Non-Gaussian Bayesian Tracking". *IEEE Transactions on Signal Processing* 50-2:174-188, February (2002).
2. Cappé, O., Moulines E. and Rydén T., *Inference in Hidden Markov Models*. Springer-Verlag (2005).
3. Desbouvieres, F. and Ait-El-Fquih, B., "Direct, prediction-based and smoothing-based particle filter algorithms". *Proc. 4th world conference of the Int. Ass. for Statistical Computing (IASC 2008)* Yokohama, Japan, Dec. 5-8 (2008).
4. Doucet, A., Godsill, S. J. and Andrieu, C., "On sequential Monte Carlo sampling methods for Bayesian filtering". *Statistics and Computing*, 10:197-208 (2000).
5. Doucet A., de Freitas N. and Gordon N. (eds.), *Sequential Monte Carlo methods in practice*. Springer Verlag (2001).
6. Gordon, N.J., Salmond, D.J. and Smith, A.F.M., "Novel approach to non linear/non-Gaussian Bayesian state estimation". *IEE Pr.-F.* 140:107-113 (1993).
7. Saha, S., Manda, P.K., Boers, Y., Driessen, H. and Bagchi, A., "Gaussian Proposal density using moment matching in SMC methods" *Statistics and Computing*. 19-2:203-208, June (2009).
8. Smith, A.F.M. and Gelfand, A.E., "Bayesian statistics without tears: a sampling-resampling perspective". *The American Statistician*. 46-2:84-87 (1992).
9. van der Merwe, R., Doucet, A., de Freitas, N. and Wan, E., "The unscented Particle Filter". *Advances in Neural Information Processing Systems*. (2001).

EM and ICE in Hidden and Triplet Markov Models

Wojciech Pieczynski

CITI Department, Institut Telecom, Telecom Sudparis
Evry, France
Email: wojciech.pieczynski@it-sudparis.eu

Abstract: This paper addresses the problem of parameter estimation in the case of hidden data. The aim is to discuss two general iterative parameter estimation methods “Expectation-Maximization” (EM) and “Iterative Conditional Estimation” (ICE) in the context of the classical Hidden Markov Models (HMMs) and in the context of the recent Triplet Markov Models (TMMs). A very general method of TMMs identification based on ICE and copulas is also specified.

Keywords: Hidden data, parameter estimation, Expectation-Maximization, Iterative Conditional Estimation, hidden Markov models, triplet Markov models, copulas.

1 Introduction

Let $Y = (Y_1, \dots, Y_n)$ be observed data and $X = (X_1, \dots, X_n)$ hidden ones. In the whole paper, each Y_i takes its values from the set of real numbers R , and each X_i takes its values from a finite set of classes $\Omega = \{\omega_1, \dots, \omega_k\}$. Let $p_\theta(x, y)$ be the probability distribution depending on a parameter $\theta \in R^m$, and let $l_\theta(x, y) = \log[p_\theta(x, y)]$ be the log-likelihood. Besides, let $\hat{\theta}(X, Y)$ be an estimator of $\theta \in R^m$ defined from complete data (X, Y) . Both “Expectation-Maximization” (EM) and “Iterative Conditional Estimation” (ICE) define a sequence of parameters from the observation y . After having chosen an initial value θ^0 , the EM sequence is defined by

$$\theta^{q+1}(y) = \arg \max_{\theta} E[l_\theta(X, Y) | Y = y, \theta^q], \quad (1.1)$$

while the ICE sequence is defined by

$$\theta^{q+1}(y) = E[\hat{\theta}(X, Y) | Y = y, \theta^q]. \quad (1.2)$$

The EM method (McLachlan and Krishnan (1997)) is well known and widely used, while ICE is less popular. However, ICE has been successfully used in different problems of unsupervised image processing; let us mention (Cao et al. (2005), Carincotte et al. (2006), Derrode and Pieczynski (2004), Destrempe and Mignotte (2004), Provost et al. (2004), and Salzenstein et al. (2007)), among recent

references. Concerning general considerations to compare EM and ICE, let us underline the following points:

(i) ICE is more general than EM because the estimator $\hat{\theta}(X, Y)$ can be of any form; in particular, it can be the “maximum likelihood” (ML) estimator or not. It is also often easier to perform because the maximization step does not exist in ICE ;

(ii) as stated in Delmas (1997), in the case of an exponential family of distributions EM and ICE can produce the same sequence (θ^q) ;

(iii) many comparisons between EM and ICE have been performed in classical contexts with Gaussian noise, like adaptive estimations (Peng and Pieczynski (1995)), hidden Markov chains (Benmiloud and Pieczynski (1995)), or hidden Markov trees (Monfrini (2002)). In all these situations the EM formulae are computable and it turns out that both EM and ICE methods are of quite a comparable efficiency ;

(iv) the use of EM is justified by the theoretical results concerning the optimal asymptotic behavior of the ML estimator, and by the fact that EM produces a sequence (θ^q) such that the sequence $p(y|\theta^q)$, being increasing, often converges to a local maximum. We have to notice that this does not imply the convergence of (θ^q) to the real parameter θ ; however, if the initial value θ^0 is close enough to the real value θ , the convergence can be shown under some mild hypotheses. The idea behind ICE is different and is based on the following. Assuming that $\hat{\theta}(X, Y)$ has interesting quadratic error - or is even optimal, being, for example, an ML estimator in an exponential model - one wishes to approximate it by a function of the only observed variables y . The “best” - with regard to the same “quadratic error” criterion - approximation is the conditional expectation. As this expectation depends on the parameter, we arrive at (1.2). Concerning the convergence of ICE, let us mention a recent theoretical result obtained in the case of independent data (Pieczynski (2008)). As in the case of EM, convergence can be obtained under some reasonable hypotheses if the initial value θ^0 is close enough to the real value θ ;

(v) EM encounters more difficulties in hidden Markov field models, where the maximization step cannot be calculated and one is obliged to simplify the model, for example by introducing the “mean field” as indicated in (Celeux et al. (2006)). ICE can be used without model modification, even in more complex situations, as in the context of recent triplet Markov fields (Benboudjema and Pieczynski (2007)).

The aim of the paper is to discuss and compare the difficulties when applying these two methods in the context of the classical Hidden Markov Models (HMMs (Cappe et al. (2005), Ephraim (2002), (Koski (2001)) and the recent Triplet Markov Models (TMMs (Pieczynski and Desbouvries (2005), (Pieczynski (2007), Pieczynski (2010)). A very general method of TMM identification based on ICE and copulas is also briefly described.

2 Pairwise and Hidden Markov Models

Let us consider the couple of stochastic sequences $(X, Y) = (X_1, Y_1, \dots, X_n, Y_n)$ and let us set $Z = (X, Y) = (Z_1, \dots, Z_n)$, with $Z_1 = (X_1, Y_1)$, ..., $Z_n = (X_n, Y_n)$. The couple $Z = (X, Y)$ is a “Pairwise Markov Model” (PMM) if its distribution is given by

$$p(z) = p(z_1)p(z_2|z_1)\dots p(z_n|z_{n-1}). \quad (2.1)$$

We will say that a PMM $Z = (X, Y)$ is “stationary” if the distributions $p(z_i, z_{i+1})$ do not depend on $i = 1, \dots, n-1$. Thus the distribution of a stationary PMM (SPMM) is given by $p(z_1, z_2)$, which can be written:

$$p(z_1, z_2) = p(x_1, x_2)p(y_1, y_2|x_1, x_2). \quad (2.2)$$

There are then two kinds of SPMM $Z = (X, Y)$. Either X is a Markov chain or it is not. If it is, the SPMM $Z = (X, Y)$ will be called a “stationary hidden Markov model” (SHMM), which is consistent with the fact that the hidden model is a Markov one. One can then show that a “reversible” (which means that $p(z_i, z_{i+1}) = p(z_{i+1}, z_i)$) SPMM is an SHMM if, and only if, $p(y_1, y_2|x_1, x_2)$ in (2.2) verifies

$$p(y_1|x_1, x_2) = p(y_1|x_1). \quad (2.3)$$

In fact, a reversible SPMM $Z = (X, Y)$ is an SHMM if, and only if, the two equivalent conditions: (i) for each $2 \leq i \leq n$, $p(y_i|x_i, x_{i-1}) = p(y_i|x_i)$; (ii) for each $1 \leq i \leq n$, $p(y_i|x) = p(y_i|x_i)$, are verified (Pieczynski (2007)).

Let us remark that the very classical SHMM, whose distribution is defined by

$$p(x, y) = p(x_1)p(x_2|x_1)\dots p(x_n|x_{n-1})p(y_1|x_1)\dots p(y_n|x_n), \quad (2.4)$$

is obtained when $p(y_1, y_2|x_1, x_2)$ in (2.2) verifies

$$p(y_1, y_2|x_1, x_2) = p(y_1|x_1)p(y_2|x_2), \quad (2.5)$$

which is stronger than (2.3).

In a similar way to the classical HMMs, the transitions $p(x_{i+1}|x_i, y)$ and the marginal distributions $p(x_i|y)$ can be computed in the following way. Let us consider the following "forward" $\alpha(x_i) = p(y_1, \dots, y_{i-1}, z_i)$ and "backward" $\beta(x_i) = p(y_{i+1}, \dots, y_n | z_i)$ probabilities, which again give the classical probabilities when the PMM considered is an HMM. Then we have

$$\alpha_1(x_1) = p(z_1), \text{ and } \alpha_{i+1}(x_{i+1}) = \sum_{x_i \in \Omega} \alpha_i(x_i) p(z_{i+1} | z_i) \text{ for } 2 \leq i \leq n; \quad (2.6)$$

$$\beta_n(x_n) = 1, \text{ and } \beta_i(x_i) = \sum_{x_{i+1} \in \Omega} \beta_{i+1}(x_{i+1}) p(z_{i+1} | z_i) \text{ for } 1 \leq i \leq n-1; \quad (2.7)$$

$$p(x_{i+1}|x_i, y) = \frac{p(z_{i+1}|z_i) \beta_{i+1}(x_{i+1})}{\beta_i(x_i)}; \quad (2.8)$$

$$p(x_i|y) = \frac{\alpha_i(x_i) \beta_i(x_i)}{\sum_{x_i' \in \Omega} \alpha_i(x_i') \beta_i(x_i')}; \quad (2.9)$$

$$p(x_i, x_{i+1}|y) = p(x_i|y) p(x_{i+1}|x_i, y). \quad (2.10)$$

The formulae (2.6)-(2.10) are extensions of the well known HMM formulae, which are obtained by taking $p(z_i) = p(x_1) p(y_1 | x_1)$ and $p(z_{i+1} | z_i) = p(x_{i+1} | x_i) p(y_{i+1} | x_{i+1})$.

Let us underline the fact that considering SPMMs which are not SHMMs (in which (2.3) does not hold) can be of real interest in the unsupervised segmentation of real or simulated data: see different results presented in (Derrode and Pieczynski (2004)).

3. EM and ICE in SPMM

Let us consider an SPMM whose distribution given by (2.2) is such that $p(y_1, y_2 | x_1, x_2)$ are Gaussian. The parameters to be estimated are $p_{jk} = p(x_1 = \omega_j, x_2 = \omega_k)$ and the mean vectors M_{jk} and variance-covariance matrices Γ_{jk} of the Gaussian distributions $p(y_1, y_2 | x_1 = \omega_j, x_2 = \omega_k)$. In both EM and ICE methods (p_{jk}) are re-estimated by

$$p_{jk}^{q+1} = \frac{1}{n-1} \sum_{i=1}^{n-1} p^q(x_i = \omega_j, x_{i+1} = \omega_k | y), \quad (3.1)$$

where $p^q(x_i = \omega_j, x_{i+1} = \omega_k | y)$ are computed with (2.6)-(2.10).

The parameters M_{jk} and Γ_{jk} are re-estimated in EM by

$$M_{jk}^{q+1} = \frac{\sum_{i=1}^{n-1} \begin{bmatrix} y_i \\ y_{i+1} \end{bmatrix} p^q(x_i = \omega_j, x_{i+1} = \omega_k | y)}{\sum_{i=1}^{n-1} p^q(x_i = \omega_j, x_{i+1} = \omega_k | y)}; \quad (3.2)$$

$$\Gamma_{jk}^{q+1} = \frac{\sum_{i=1}^{n-1} \left(\begin{bmatrix} y_i \\ y_{i+1} \end{bmatrix} - M_{jk}^{q+1} \right) \left(\begin{bmatrix} y_i \\ y_{i+1} \end{bmatrix} - M_{jk}^{q+1} \right)^T p^q(x_i = \omega_j, x_{i+1} = \omega_k | y)}{\sum_{i=1}^{n-1} p^q(x_i = \omega_j, x_{i+1} = \omega_k | y)}, \quad (3.3)$$

while in ICE they are re-estimated by

$$M_{jk}^{q+1} = \frac{\sum_{i=1}^{n-1} \begin{bmatrix} y_i \\ y_{i+1} \end{bmatrix} 1_{[x_i^q = \omega_j, x_{i+1}^q = \omega_k]}}{\sum_{i=1}^{n-1} 1_{[x_i^q = \omega_j, x_{i+1}^q = \omega_k]}}; \quad (3.4)$$

$$\Gamma_{jk}^{q+1} = \frac{\sum_{i=1}^{n-1} \left(\begin{bmatrix} y_i \\ y_{i+1} \end{bmatrix} - M_{jk}^{q+1} \right) \left(\begin{bmatrix} y_i \\ y_{i+1} \end{bmatrix} - M_{jk}^{q+1} \right)^T 1_{[x_i^q = \omega_j, x_{i+1}^q = \omega_k]}}{\sum_{i=1}^{n-1} 1_{[x_i^q = \omega_j, x_{i+1}^q = \omega_k]}}, \quad (3.5)$$

where $x^q = (x_1^q, \dots, x_n^q)$ is sampled according to $p(x|y)$ using the current values of the parameters.

Dealing with the Gaussian case under consideration here with either EM or ICE would probably provide similar results, as they do in the classical HMMs. However, when one leaves the Gaussian case and deals with the “generalized” mixture estimation, ICE is much easier to apply. In the SHMM context one is faced with the “generalized” mixture estimation problem when the forms of the noise distributions $p(y_1 | x_1 = \omega_j)$ are not known and can vary with the class ω_j . However, for each ω_j one knows that $p(y_1 | x_1 = \omega_j)$ belongs to a given set of forms. For example, one knows that $p(y_1 | x_1 = \omega_1)$ is either Gaussian or gamma, $p(y_1 | x_1 = \omega_2)$ can be Gaussian, exponential, or Rayleigh, ... and so on. Such situations are of interest and they can occur, in particular, in radar images models (Delignon and Pieczynski (2002), Nadarajah and Kotz (2008)). Estimating such a

mixture therefore contains two problems: (i) finding the right form for each class; and (ii) estimating the related parameters. ICE has been extended to a “generalized” ICE (GICE) to deal with such problems in SHMMs in (Giordana and Pieczynski (1997)) and different experiments have shown its efficiency. Afterwards, the extension of ICE to “generalized” reversible SPMMs has been suggested in (Pieczynski (2010)). Let us briefly recall its principle.

For K classes there are $(K-1)K/2$ distributions $p(y_1, y_2 | x_1, x_2)$ on R^2 . Besides, let $H(y_1, y_2)$ be a cumulative distribution function (cdf) over R^2 , and $H_1(y_1)$, $H_2(y_2)$ the related marginal cdfs. Then, according to the Sklar theorem (Brunel and Pieczynski (2005), Nelsen (1998)), there is a unique cdf C on $[0,1]^2$ with uniform marginal distributions (called a “copula”) such that

$$H(y_1, y_2) = C(H_1(y_1), H_2(y_2)) \quad (3.6)$$

Thus each of the $(K-1)K/2$ distributions $p^{ij}(y_1, y_2) = p(y_1, y_2 | x_1 = \omega_i, x_2 = \omega_j)$ on R^2 is defined by $(K-1)K/2$ marginal distributions $p^{ij}(y_1)$ and $(K-1)K/2$ copulas C^{ij} . Assuming that for each (i, j) the form of the marginal distribution $p^{ij}(y_1)$ belongs to a given set $\Phi^{ij} = \{F_1^{ij}, \dots, F_{r(i,j)}^{ij}\}$ of admissible forms and the form of the copula C^{ij} belongs to a given set $X^{ij} = \{C_1^{ij}, \dots, C_{m(i,j)}^{ij}\}$ of admissible forms, one is faced with the following problem : for each (i, j) select from Φ^{ij} and X^{ij} the correct forms and estimate the related parameters. At each iteration of ICE these two problems are then dealt with using $x^q = (x_1^q, \dots, x_n^q)$ sampled according to $p(x|y, \theta^q)$.

4. Generalized ICE in Stationary Triplet Markov Models

Let us consider the couple (X, Y) as above. Let $U = (U_1, \dots, U_n)$ be a third random chain, each U_i taking its values from $\Lambda = \{\lambda_1, \dots, \lambda_M\}$. The triplet $T = (X, U, Y)$ is called a “Triplet Markov Model” (TMM) if its distribution is a Markovian one. Setting $V = (X, U)$ one sees that a TMM can also be seen as a PMM (V, Y) ; in fact, $V = (V_1, \dots, V_n)$ with each V_i taking its values from a finite set $\Omega \times \Lambda$. Thus both X and U can be estimated by some Bayesian method, and the parameters can be estimated with EM or ICE as discussed above.

The choice of the interpretation of the third chain U and the choice of the Markovian distribution for $T = (X, U, Y)$ lead to a very rich family of possible distributions for (X, Y) . One possible choice is a Markov distribution for $V = (X, U)$ such that X is a semi-Markov chain (Pieczynski and Desbouvries (2005)); $T = (X, U, Y)$ is then a classical hidden semi-Markov chain. Other

choices lead to a non-stationary distribution for (X, Y) , where the switches among the different stationarities are modeled by U (Lanchantin and Pieczynski (2004)). Let us also mention the use of TMM to perform the Dempster-Shafer fusion in a Markovian context (Pieczynski (2007)). It is also possible to consider multivariate U , to model different properties simultaneously. For example, one can take $U = (U^1, U^2)$, where U^1 models the semi-Markovianity of X and U^2 models its non-stationarity (Lapuyade-Lahorgue and Pieczynski (2006)). In each of these situations, one can then apply the “generalized” ICE described above to the related PMM $T = (X, U, Y) = (V, Y)$.

References

1. Benboudjema D., and Pieczynski, W., Unsupervised statistical segmentation of non stationary images using triplet Markov fields, *IEEE Trans. on Pattern Analysis and Machine Intelligence*, **29**, 8, 1367-1378 (2007).
2. Benmiloud, B. and Pieczynski, W., Estimation des paramètres dans les chaînes de Markov cachées et segmentation d'images, *Traitement du Signal*, **12**, 5, 433-454 (1995).
3. Brunel, N., and Pieczynski, W., Unsupervised signal restoration using hidden Markov chains with copulas, *Signal Processing*, **85**, 12, 2304-2315 (2005).
4. Cao, Y. F., Sun, H., and Xu, X., An unsupervised segmentation method based on MPM for SAR images, *IEEE Geoscience and Remote Sensing Letters*, **2**, 1, 55-58 (2005).
5. Cappé, O., Moulines, E., and Ryden, T., *Inference in hidden Markov models*, Springer, Series in Statistics, Springer (2005).
6. Carincotte, C., Derrode, S., and Bourennane, S., Unsupervised change detection on SAR images using fuzzy hidden Markov chains, *IEEE Trans. on Geoscience and Remote Sensing*, **44**, 2, 432-441 (2006).
7. Celeux, G., Forbes, F., and Peyrard, N., EM procedures using mean field-like approximations for Markov model-based segmentation, *Pattern Recognition*, **36**, 131-144 (2006).
8. Delignon, Y., and Pieczynski, W., Modeling non Rayleigh speckle distribution in SAR images, *IEEE Trans. on Geoscience and Remote Sensing*, **40**, 6, 1430-1435 (2002).
9. Delmas, J.-P., An equivalence of the EM and ICE algorithm for exponential family, *IEEE Trans. on Signal Processing*, **45**, 10, 2613-2615 (1997).
10. Derrode, S., and Pieczynski, W., Signal and image segmentation using pairwise Markov chains, *IEEE Trans. on Signal Processing*, **52**, 9, 2477-2489 (2004).
11. Destrempe, F., and Mignotte, M., A statistical model for contours in images, *IEEE Trans. on Pattern Analysis and Machine Intelligence*, **26**, 5, 626-638 (2004).

12. Ephraim, Y., Hidden Markov processes, *IEEE Trans. on Information Theory*, **48**, 6, 1518-1569 (2002).
13. Giordana, N., and Pieczynski, W., Estimation of generalized multisensor hidden Markov chains and unsupervised image segmentation, *IEEE Trans. on Pattern Analysis and Machine Intelligence*, **19**, 5, 465-475 (1997).
14. Koski, T., *Hidden Markov models for bioinformatics*, Kluwer Academic Publishers, Netherlands (2001).
15. Lanchantin, P., and Pieczynski, W., Unsupervised non stationary image segmentation using triplet Markov chains, ACVIS 04, Aug. 31-Sept. 3, Brussels, Belgium, 2004.
16. Lapuyade-Lahorgue, J., and Pieczynski, W., Unsupervised segmentation of hidden semi-Markov non stationary chains, MaxEnt 2006, Paris, France, July 8-13, 2006.
17. McLachlan, G. J. and Krishnan, T., *EM Algorithm and Extensions*, Wiley (1997).
18. Monfrini, E., Identifiabilité et méthode des moments dans les mélanges généralisés de distributions du système de Pearson, Thèse de l'Université Paris VI, soutenue le 4 janvier 2002.
19. Nadarajah, S., and Kotz, S., Intensity models for non-Rayleigh speckle distributions, *International Journal of Remote Sensing*, **29**, 2, 529-541 (2008).
20. Nelsen, R. B., *An introduction to Copulas*. Number 139 in Lecture Notes in Statistics. Springer-Verlag (1998).
21. Peng, A. and Pieczynski, W., Adaptive mixture estimation and unsupervised local Bayesian image segmentation, *Graphical Models and Image Processing*, **57**, 5, 389-399 (1995).
22. Pieczynski, W., and Desbouvries, F., On triplet Markov chains, (ASMDA 2005), Brest, France, May 2005.
23. Pieczynski, W., Multisensor triplet Markov chains and theory of evidence, *International Journal of Approximate Reasoning*, **45**, 1, 1-16 (2007).
24. Pieczynski, W., Sur la convergence de l'estimation conditionnelle itérative, *Comptes Rendus*, **346**, 7-8, 457-460 (2008).
25. Pieczynski, W., *Triplet Markov chains and image segmentation*, chapter 4 in *Inverse Problems in Vision and 3D Tomography*, A. Mohammed-Djafari ed., Wiley (2010).
26. Provost, J.-N., Collet, C., Rostaing, P., Pérez, P., and Bouthemy, P., Hierarchical Markovian segmentation of multispectral images for reconstruction of water depth maps, *Computer Vision and Image Understanding*, **93**, 2, 155-174 (2004).
27. Salzenstein F., Collet, C., Le Cam, S. and Hatt, M., Non stationary fuzzy Markov chains, *Pattern Recognition Letters*, **28**, 16, 2201-2208 (2007).

High school schedule creation, optimization problems and solution

Lina Pupeikienė

Vilnius Gediminas technical University
Saulėtekio al. 11, LT-10223 Vilnius, Lithuania,
Lina.Pupeikiene@gmail.com

Abstract. The timetable problem, searching the timetable for the class assignment in the schools, belongs to the group of NP hard and NP complete problems. In Lithuanian high schools every pupil can choose a lot of subjects by his wish. The problem is more complicated when the every pupil has possibilities to choose not only subjects, but hour per week of this subject too. However, as the number of teachers, number of pupils, number of different subjects, number of different subject hours, time slots, and the constraints increases, the required time to find at least one feasible solution grows exponentially. Global optimization algorithms are a quite common approach to solve this problem. In this paper, we describe the advantages of distributed school schedule optimization a software system was developed using Java technology and grid computing techniques. Optimization algorithms used in the software include the Monte-Carlo local optimization algorithm, the Simulated Annealing and Bayes global optimization algorithms.

Key words: Global optimization, school schedule creation, distributed schedule optimization, Monte-Carlo, Simulated Annealing, Bayes.

1 Introduction

A timetable specifies which people meet at which location and at what time. The timing of events must be such that nobody has more than one event at the same time. School timetabling as a term refers to the construction of weekly timetables for schools of secondary education [14]. Specific feature of school timetabling field is a great number of research papers and widely used commercial software. Therefore a discussion of new results will be.

The events are lessons in a subject, taught by a teacher to a group of pupils in a single room. The timetable assigns a teacher, a pupils group, a room, and a time slot to each lesson. The pupil groups are specific to the subject, we call them subject-groups. A high school is referred here as the last grades of a high school or gymnasium where the pupils can mostly choose their preferred learning profile subjects. Therefore, this task is more complex in comparison with a secondary school scheduling without high school classes.

Some combinations of assignments lead to acceptable timetables, constraints follow from conditions imposed by rooms, pupils or teachers. We distinguish two types of constraints: conditions that must be met ("hard" constraints) and desires that should be fulfilled as well as possible ("soft" constraints). An important set of soft constraints is defined by didactic reasons. For example, by placing "hard" subjects, such as mathematics or physics, into morning hours. The maximal number of daily hours T_{max} is obviously a hard constraint. Timetabling can be generally defined as the activity of assigning, subject to constraints, a number of events to a limited number of time periods and locations such, that desirable

objectives are satisfied as nearly as possible [26]. Educational timetabling can be divided into three main classes: school timetabling, course timetabling and exam timetabling [15]. The goal is to find a timetable that satisfies all the hard constraints and minimizes the violation of soft constraints.

2 Overview of publications

A survey on educational timetabling problems [23] gives an overview of the literature. Overviews on examination timetabling and university course timetabling are in [4, 12, 13]. A comprehensive overview of formulations and of state-of-the-art approaches is in the surveys [4, 7, 8, 13, 15], in the proceedings of the PATAT conferences [5 – 7, 9, 10] and in the Lecture Notes in Computer Science series [9 – 11]. The European working group on automated timetabling (EURO-WATT) maintains a website with information on timetabling problems [25].

3 New Elements

The first new element of this work is the application and systematic investigation of the Bayesian Heuristic Approach [20] for optimization of heuristic parameters. These include the initial temperature and the cooling rate of Simulating Annealing (SA) algorithm and the randomization parameter of the local search algorithm. The formulation of the objective function in terms of Pareto optimality seems to be new in the field of school scheduling. The paper describes apparently the first web-based platform-independent implementation of the software. Java servlet provides conditions for application at any school with internet connection. Any web browser works, no additional software is needed. Note that efficiency of recent versions of Java is close to that of the most efficient programming languages [9].

4 Defining Optimization Problem

Ministry of Education of the Republic of Lithuania has confirmed basic rules for high school schedule forming. They can be complementary of each school's rules and restrictions. However, the main purpose of these limitations is to develop a schedule, which would evaluate of the Ministry of Education requirements. In addition, this schedule must be acceptable to both: pupils and teachers.

Required schedule restrictions (formed by the Ministry of Education):

- * Working days d per week must be $d \leq 5$.
- * The teacher simultaneously cannot work in several different places.
- * The teacher cannot have more than 36 hours per week.
- * The pupil simultaneously cannot learn few different subjects.
- * A pupil i may have $28 \leq i \leq 32$ lessons per week.
- * It cannot be more than $p \leq 7$ lessons p per day.
- * Number of pupils i in one subject-group can be $15 \leq i \leq 30$.
- * In each classroom simultaneously cannot be several different types of subjects (for example, mathematics and physics).

* Subjects, requiring special measures or facilities, shall be taught in the special classrooms (for example, IT, chemistry etc.).

Technically any required restriction violations cannot be broken. There can be only some minor offenses necessary restrictions, if it significantly improves the quality of the schedule. To define with timetable is good or bad we use penalty points. The penalty point's c_r , which assessing these restrictions, should be imposed very strictly.

The main required penalty point's restrictions function is as follows:

$$F_f = \sum_r c_r N_r$$

here c_r – penalty for required restriction r ; N_r – number of required restriction. In this case $r = 1, \dots, 9$.

Some of required restrictions c_r can be evaluated by the individual rules of each school. Such requirements are called needful, or “soft” constrains. They are valued differently in each school.

The main needful restrictions of the schedule include:

- * Elimination of “windows” in teacher’s schedule.
- * Elimination of “windows” in pupil’s schedule.
- * Unacceptable working hours.
- * Unacceptable workdays.
- * Unacceptable order of subjects.
- * Changing of pupils in the formed subject-group.

Usually penalty points for these restrictions are as follows:

c_m – penalty for the “window” on teacher’s m schedule.

c_s – penalty for the “window” on pupils s schedule.

c_{mv} – penalty for “bad” hour v of teacher m .

c_{md} – penalty for “bad” day d of teacher m .

c_{sv} – penalty for “bad” hour v of pupil s .

c_{pd} – penalty for violation of pedagogical didactic pd .

c_{mg} – penalty of the list change of subject-group g taught by teacher m .

“Bad” hour/day is the hour/day, when teacher/pupil already has a work hour.

Pedagogical didactic evaluates the difficulty of subjects. Most difficult subjects must be in the 1-4 lessons during the day. Less important subjects – in the end of the day. The importance of every subject is written in initial data file.

The sum function of the needed restrictions penalty points is as follows:

$$F_n = \sum_m c_m L_m + \sum_s c_s L_s + \sum_m \sum_v c_{mv} L_m^v + \\ + \sum_m \sum_d c_{md} L_m^d + \sum_s \sum_v c_{sv} L_s^v + \sum_{pd} c_{pd} L_{pd} + \sum_n c_{ng} L_n,$$

here L_m – number of “windows” on teachers m schedule; L_s – number of “windows” on pupils s schedule; L_m^v – number of “bad” hours v on the teachers m schedule; L_m^d – number of “bad” days d on the teachers m schedule; L_s^v – number

of “bad” hours v on the pupils s schedule; L_{pd} – number of pedagogical didactic pd violations; L_n – number n of changing formed subject-group.
 All physical restrictions and inconveniences are showed in Figure 1.

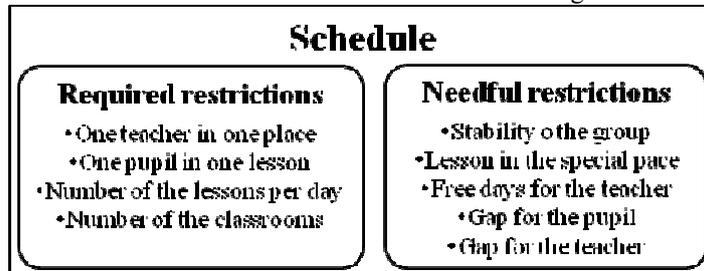


Figure 1. Restrictions for a creation of high school schedule

A compromise solution is reached by defining penalties for violation of constraints and disregarding inconveniences. Therefore, penalty points are calculated:

$$F = F_f + F_n$$

where, F_f – is a sum of the penalties for the required restrictions; F_n – is a sum of the penalties for the needful restrictions (disregarding inconveniences). Optimal schedule will be schedule, which has as less as possible penalty points. To find such schedule, objective function F should be optimized. To not analyze the schedules with same number of penalty points, Pareto optimality was formulated. So we will get less variants to analyze and will save the users time. The optimization problem is

$$\min_{\tau \in A} F(\tau)$$

where, $F(\tau)$ is the total penalty of some schedule τ ; A is the set of schedules satisfying the physical constraints. The penalties $F(\tau)$ depend on expert evaluations, therefore we regard them as heuristics.

5. Optimization Methods

5.1 Defining Neighborhood

Many different definitions can be used defining neighbourhood in a set A of feasible timetables d . The definition is important because local search is performed in the neighbourhood of the given point. We search for better timetables by subsequent closing of gaps for pupils and teachers. In this case the neighbours of a timetable d' are all timetables d'' that can be reached from d' by a sequence of closing gap operations. This way we obtain locally optimal $d^*(d')$ that depends on the initial point d' .

Local search can be randomized by selecting current candidate (a pupil or a teacher) for gap closing with some probability x_0 . Closing gaps for randomly selected pupils and teachers, we modify the search sequences. However, this not helps to reach the global optimum since the neighbourhood remains the same.

5.2 Escaping Neighborhood

Simplest algorithm to search for global optimum is just random search with uniform distribution of observations (observation is calculation of the objective function at some fixed point). The advantages are simplicity and convergence to a global minimum of continuous functions. A well-known way to escape the local minimum is Simulated Annealing [1, 2, 14, 19, 21, 22]. Denote

$$\delta_n = F(d^{n+1}) - F(d^n)$$

Here d^n is a current timetable, d^{n+1} is a new timetable generated by closing gap operation. Define the probability

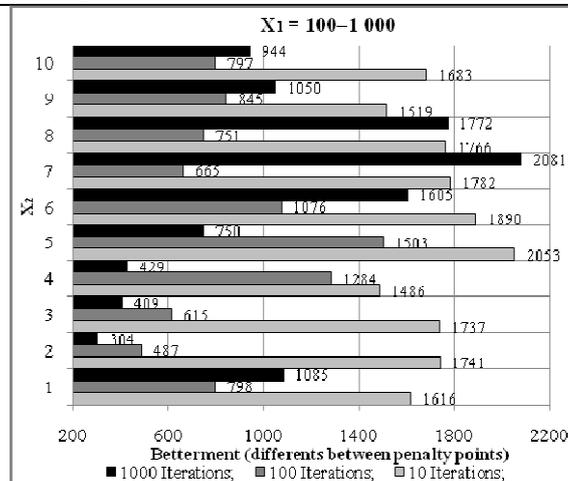
$$p_n = e^{-\frac{\delta_n}{x_1 / \ln(1+x_2^n)}}, \text{ if } \delta_n > 0,$$

$$p_n = 1, \text{ if } \delta_n < 0,$$

where parameter x_1 is the “initial temperature”, parameter x_2 defines the “cooling rate”. SA algorithm means:

go to new timetable d^{n+1} with probability p_n

To apply the SA to a specific problem, one must specify the parameters x_1 and x_2 . The choice can have a significant impact on the method's effectiveness. Unfortunately, there are no choices of these parameters that will be good for all problems. Analyzing Figure 2, we see different results using different initial parameters. Here difference of penalty points (between initial and optimal schedules) is calculated. Every column is received after 100 experiments with fixed initial parameters (Iterations, x_1 and x_2). In the left side of Figure 2 the results are grouped by x_2 when x_1 was between 100 and 1000. In the right side, the results are grouped by x_1 when x_2 was between 1 and 10. There are showed only best results.



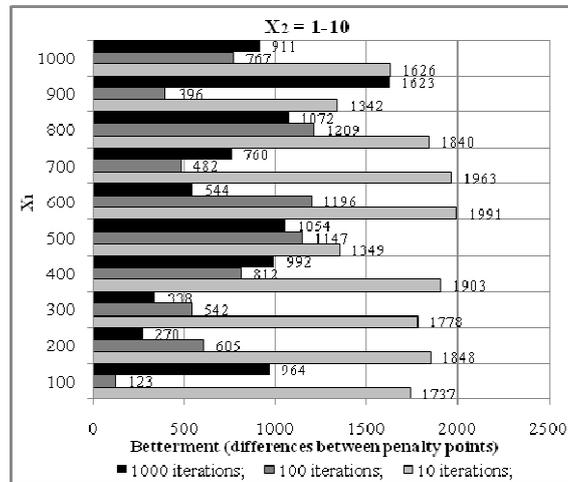


Figure 2. The best results of SA using different parameters

We cannot see optimal parameters x_1, x_2 of SA. Optimal results depend on the initial soft constrains and number of iteration. A way to adapt these parameters to a given problem is automatic optimization. This is not an easy problem since we need optimize multi-modal function with considerable noise. Here the Bayesian Heuristic Approach (BHA) [20] is useful. Figures 3 and 4 illustrate efficiency of automatic adaptation of SA parameters using BHA. In these figures, the difference between initial and optimal timetable is showed. There we see 100 experiments with every different SA iteration. SA parameters were set automatically. Figure 3 shows, that method is more efficient as more SA iteration are used. Figure 4 illustrates the best results what was shown during 100 experiments with every different SA iteration. There we can see, that the best results we will get when it will be many SA and BHA iterations.

5.3 Bayesian Heuristic Approach

The Bayesian Heuristic Approach was designed for automatic optimization of heuristic parameters by filtering the noise during optimization of multi-modal functions [20]. We need to optimize three heuristic parameters $x = (x_0, x_1, x_2)$. Optimal parameters are obtained using the data of some specific school.

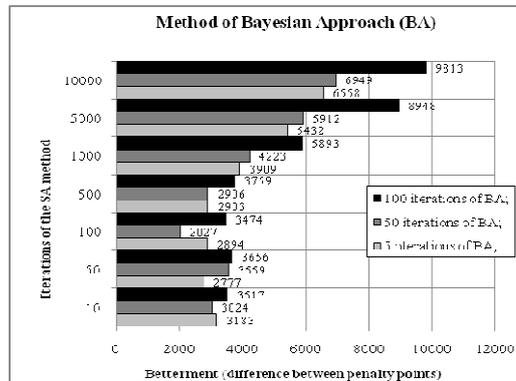


Figure 3. Average of 100 experiments results using BHA

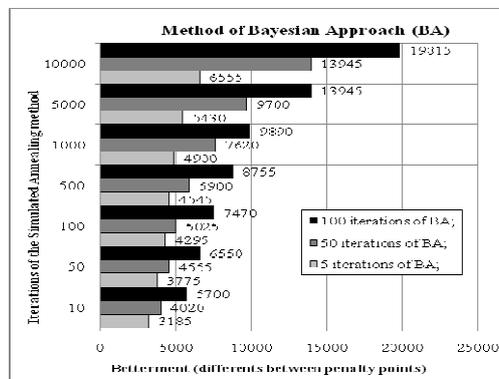


Figure 4. The best results of BHA after 100 experiments with each different SA iteration

However, the results can be used in similar schools as an approximation.

6 “School schedule optimization” program working steps

“School schedule optimization” program designed to high school scheduling.

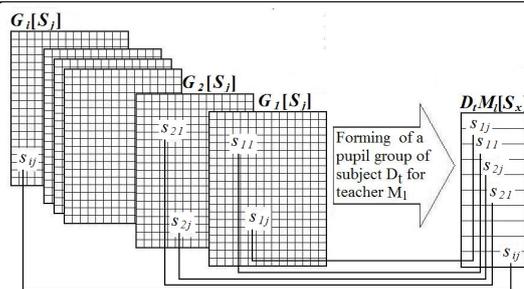


Figure 5. Forming subject-groups to teachers

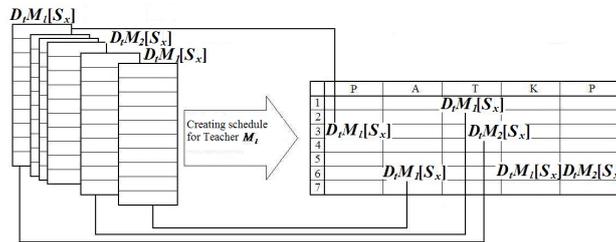


Figure 6. Time table for teachers creation

Figure 5 illustrates how subject-groups are assigned to teachers. Here pupils s_{ij} from groups G_i are grouped to the groups with identical subject D_i . Identical subject has same name and same hours per week. These groups are called subject-groups (with x pupils in the group) and assigned to the teacher M_i . Figure 6 shows how teacher's timetables are created. The subject-groups $D_i M_i [S_x]$, with teacher M_i , subject D_i and pupils of this subject-group S_x , are putted to the free class-room and to school timetable. When process is finalized, the optimization process is ready to start.

After optimizing, we can see such results of this program:

- * school schedule;
- * individual pupils schedules;
- * individual teachers schedules;
- * individual room schedules;
- * subject-group schedules;

All results user can see in the program (on working time), or download them as archive personal computer. The program does not require much effort to the user, the payment to work with a computer, or a lot of time to understand how system works.

7 Comparison of results

Here are compared such results: real schedule created in a Lithuanian high school and, from pupils and teachers wishes, created and optimized schedule. Schedule was automatically optimized with Bayes method. The results we can see in Figure 7. Both, schedule and data are from the same school and same classes. Evaluating both types of schedules, penalty points were calculating for:

- * pupil window – 5;
- * teacher window – 300;
- * teachers wished free time – 10;
- * exceeding maximum hour limit – 2000;
- * pedagogical didactic – 5.

Sum of seted penalty points for the real schedule was 380 020. It is always same, because after finishing the creation process it can't be changed. Sum of penalty points after optimization process (was seted same penalty points) are showed in the Figure 7. There are few results after optimization with different initial parameters of optimisation method Bayes. The results are different while every time schedule is created from the new point.

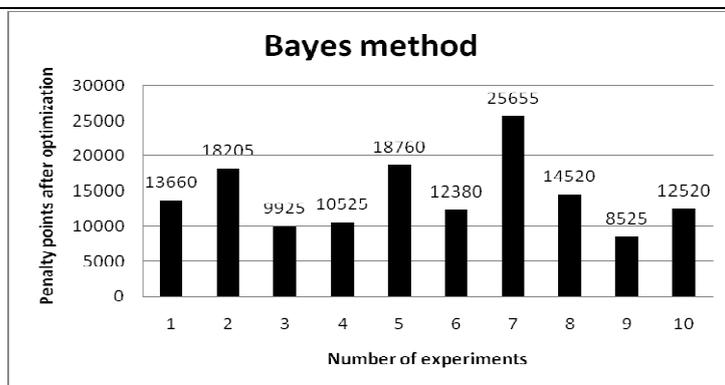


Figure 7. Penalty points after creating and optimizing schedule from initial data file

As we can see, the optimization results are much better as real schedule result. It is so, while optimization program creates and optimizes schedule only for high school classes. However, teacher can work in basic school to. However, in Lithuanian schools schedule creating starts from high school classes schedule. “School schedule optimization” program is working same way.

8 Optimization in Commercial Software

We discuss optimization possibilities of the following three commercial timetabling systems currently used in Lithuanian high schools: “Mimosa 2009”, “aSc TimeTables 2009”, and “Rector 2009”. “Mimosa 2009” [18] is the product of the Finnish company “Mimosa Software Ltd”. “Mimosa” provides convenient GUI for manual timetabling and reports constraints violations. Figure 8 shows a fragment of the output. In the upper-left side we can see pupils schedule, under it – pupils of the subject-group and in the right side – individual schedules of every pupil in the subject-group. The form is acceptable for Lithuanian schools. For example, “Ch3BK” means a chemistry lessons, pupils from 3-rd level, will learn as basic course.

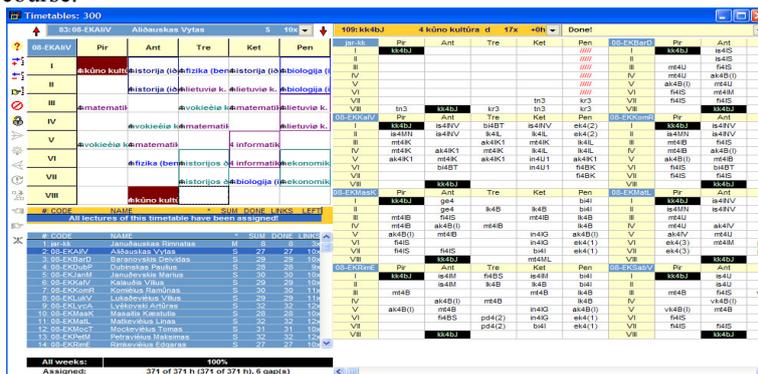


Figure 8. A fragment of “Mimosa 2009” output

Optimization is limited to closing some gaps in teacher’s schedules. The software is popular in basic schools. Application in upper classes of high schools is possible within some strict limitations by setting individual pupil schedules. Long and hard manual work is needed if the school is large. Any penalty points are calculated in this program.

“Rector 2009” [24] is the product of the Russian company “P. Yu. Smykalov”. Figure 9 shows a fragment of output in the format similar to MS “Excel” forms used in local schools. In the upper side the subject for the group 12a are showed. Under it – all groups, lessons per week, subjects and teachers are showed. Green colour means, that no one works at the same time in two places. Reports, if one is trying to insert data to wrong place, are showed in red colour. Convenient for basic school scheduling. No automatic optimization.

The screenshot shows the 'Rector 2009' software interface. At the top, there are tabs for 'Sarakaba', 'Kārtas', 'Tvaikarāšis', and 'Pakelme'. Below that, there are navigation buttons and a dropdown menu showing '12a'. The main area displays a subject schedule for group 12a, with columns for 'Pirmādiens', 'Antrādiens', 'Trešādiens', 'Ceturtdiēns', and 'Piektdiēns'. The subjects listed are: 1. Ietuvu, 2. Bioloģija, 3. Informātika, 4. Ģeogrāfija, 5. Istorija, 6. Angļu, 7. 1.kurso kultūra, 2.kurso kultūra, and 3.vokāli. Below this, there is a table with columns: 'Klasist(s)', 'Valodā', 'Nr.', 'Dzītkas', 'Mokotājs', 'Kabinets', 'Poro', 'Pamoka', 'Savate', 'Daf/VK', and 'M'. The table lists various groups (12a, 12b, 12c, 12d, 12e, 12f, 12g, 12h, 12i, 12j, 12k, 12l, 12m, 12n, 12o, 12p, 12q, 12r, 12s, 12t, 12u, 12v, 12w, 12x, 12y, 12z) and their corresponding teachers and subjects.

Figure 9. A fragment of “Rector 2009” output

The screenshot shows a colorful grid representing a time table output from 'aSc TimeTable 2009' software. The grid is organized into two main sections: 'PIRMADIENS' (Monday) and 'ANTRADIENS' (Tuesday). Each section has columns for days of the week (1, 2, 3, 4, 5, 6, 7) and rows for different groups (11a, 11b, 11c, 11d, 11e, 11f, 11g, 12a, 12b, 12c, 12d, 12e, 12f, 12g). The cells in the grid contain letters and symbols representing subjects and teachers. For example, group 11a has 'Geo' on Monday, 'Mat' on Tuesday, 'Lk' on Wednesday, 'Ist' on Thursday, 'Ist' on Friday, 'A' on Saturday, 'Pr' on Sunday, 'A' on Monday, 'Info' on Tuesday, 'Lk' on Wednesday, 'Biol' on Thursday, 'R' on Friday, 'Ist' on Saturday, and 'V' on Sunday. The colors of the cells indicate different subjects or teachers.

Figure 10. A fragment of “aSc TimeTable 2009” output

“aSc TimeTables 2009” [3] is the product of the Slovak company “Applied Software Consultants s.r.o”. A fragment of resulting timetable for Monday and Tuesday in a compact form for eight pupil subject-groups is in Figure 10. The results of experimental calculations are in Table 1. They show that the software works well in basic schools and is not practical in large high schools. Any penalty points are calculated.

Table 1. Testing „aSc TimeTables 2009“

		Small high school (50 pupils)			Medium high school (150 pupils)			Large high school (350 pupils)						
Settings of the program		Restrictions	Complexity			Restrictions	Complexity			Restrictions	Complexity			
			Soft	Average	Hard		Soft	Average	Hard		Soft	Average	Hard	
Calculations		Time	Left subjects	Viewed options			Left subjects	Viewed options			Left subjects	Viewed options		
				---	---	---		---	---	---		---	---	---
		00:08:40	31	142800	Small	Soft	---	---	---	---	---	---	---	
		00:13:02	25	183265	Average	Soft	---	---	---	---	---	---	---	
		02:10:07	29	1652362	Hard	Soft	---	---	---	---	---	---	---	
		00:41:32	32	1375124	Small	Strict	---	---	---	---	---	---	---	
		05:53:13	27	5026351	Average	Strict	---	---	---	---	---	---	---	
		71:23:51	33	70054852	Hard	Strict	---	---	---	---	---	---	---	
		---	52	975236	Small	Soft	---	---	---	---	---	---	---	
		---	45	4523625	Average	Soft	---	---	---	---	---	---	---	
		---	---	---	Hard	Soft	---	---	---	---	---	---	---	
		---	52	33245895	Small	Strict	---	---	---	---	---	---	---	
		---	---	---	Average	Strict	---	---	---	---	---	---	---	
		---	---	---	Hard	Strict	---	---	---	---	---	---	---	
		---	---	---	Small	Soft	---	---	---	---	---	---	---	
		---	---	---	Average	Soft	---	---	---	---	---	---	---	
		---	---	---	Hard	Soft	---	---	---	---	---	---	---	
		---	---	---	Small	Strict	---	---	---	---	---	---	---	
		---	---	---	Average	Strict	---	---	---	---	---	---	---	
		---	---	---	Hard	Strict	---	---	---	---	---	---	---	

A timetable that satisfies all necessary conditions is regarded as feasible. A feasible timetable is optimal if it minimizes all undesirable factors. To compare the quality of different feasible timetables we must evaluate at least the most important undesirable factors. The difficulty is that desirability is subjective by definition and depends on the local conditions. This prevents comparison of results obtained by automatic optimization with decisions made by human operator.

To compare results of different automatic optimization methods we need procedures for evaluation of undesirable factors in some fixed scales. In this paper, it is done in the framework of Pareto optimality [16]. The commercial software does not support this, since no direct comparison of decisions quality cannot be made.

9 Concluding Remarks

- * The new element of this work is application and systematic investigation of the Bayesian Heuristic Approach (BHA) [20] to optimization of heuristic parameters (with penalty points). These include the initial temperature and the cooling rate of SA algorithm and the randomization parameter of the local search algorithm.
- * BHA is intended for global optimization of functions with noise what is typical in optimization of heuristic parameters.
- * The formulation of the objective function in terms of Pareto optimality seems to be new in the field of school scheduling.
- * Application in some large schools shows some advantages comparing with commercial software. The web-site: <http://soften.ktu.lt/~mockus> and accompanying web-sites include corresponding.

References

1. Aarts, E.H.L.; Van Laarhoven, P.J.M. 1987. *Simulated annealing: Theory and applications*. D. Reidel, Dordrecht, The Netherlands.
2. Abramson, D. 1991. *Constructing school timetables using simulated annealing: Sequential and parallel algorithms*. *Management Science*, 37:98-113.
3. aSc Timetables. 2009. <<http://www.asctimetables.com/>>
4. Bardadym, V. A. 2003. *Computer-aided school and university timetabling: The new wave?* In *Lecture notes in computer science: Vol. 1153, Practice and Theory of Automated Timetabling*, First International Conference, Selected papers, pages 22-45. Springer.
5. Burke, E. K.; Trick, M. 2005. *Practice and theory of automated timetabling V*. *Lecture notes in computer science*, Vol. 3616. Springer, Berlin.
6. Burke, E. K.; Eckersley, A. J.; McCollum, B.; Petrovic, S.; Qu, R. 2004. *Analysing similarity in examination timetabling*. In *Proceedings of the 5th International Conference on the Practice and Theory of Automated Timetabling*, pages 557-559. Springer.
7. Burke, E. K.; De Causmaecker, P. 2003. *Practice and theory of automated timetabling IV*. *Lecture notes in computer science*, Vol. 2740. Springer, Berlin.
8. Burke, E. K.; Petrovic, S. 2002. *Recent research directions in automated timetabling*. *Journal of Operational Research Society*., 140:266-280.
9. Burke, E. K.; Erben, W. 2001. *Practice and theory of automated timetabling III*. *Lecture notes in computer science*, Vol. 2079. Springer, Berlin.
10. Burke, E. K.; Carter, M. W. 1998. *Practice and theory of automated timetabling II*. *Lecture notes in computer science*, Vol. 1408. Springer, Berlin.
11. Burke, E. K.; Ross, P. 1996. *Practice and theory of automated timetabling*. *Lecture notes in computer science*, Vol. 1153. Springer, Berlin.

12. Carter, M. W.; Laporte, G. 1998. *Recent developments in practical course timetabling*. In Lecture notes in computer science: Vol. 1408. Practice and theory of automated timetabling, pages 3-19. Berlin, Springer.
13. Carter, M. W.; Laporte, G. 1996. *Recent developments in practical examination timetabling*. In Lecture notes in computer science: Vol. 1153. Practice and theory of automated timetabling, pages 3-21. Springer.
14. Dascalki, S.; Birbas, T.; Housos, E. 2004. *An integer programming formulation for a case study in university timetabling*. European Journal of Operational Research, 153:117-135.
15. De Werra, D. 1985. *An introduction to timetabling*. European Journal of Operational Research, 19:151-162.
16. Fudenberg, D.; Tirole, J. 1983. *Game Theory*. MIT Press, Boston.
17. Metropolis, N.; Rosenbluth, A. W.; Rosenbluth, M. N.; Teller, A. H.; Teller, E. 1953. *Equations of state calculation by fast computing machines*. Journal of Chemical Physics, 21:1087-1092.
18. Mimosa scheduling software. 2009. <<http://www.mimosasoftware.com/>>
19. Mockus, J. 2002. *Bayesian heuristic approach to global optimization and examples*. Journal of Global Optimization, 22:191-203.
20. Mockus, J.; Eddy, W.; Mockus, A.; Mockus, L.; Reklaitis, G. 1997. *Bayesian Heuristic Approach to Discrete and Global Optimization*. Kluwer Academic Publishers, ISBN 0-7923-4327-1, Dordrecht-London-Boston.
21. Pedroso, J. P.; Moreira, N.; Reis, R. 2004. *A web-based system for multi-agent interactive timetabling*. In ICKEDS'04: International Conference on Knowledge Engineering and Decision Support, pages 1-6, Porto, Portugal.
22. Pupeikienė, L.; Mockus, J. 2005. *School schedule optimization program*. Information Technology and Control, 34:161-170.
23. Schaerf, A. 1995. *A survey of automated timetabling, technical report cs-r 9567*. Technical report, Centrum Wiskunde & Informatica (CWI), Amsterdam, The Netherlands.
24. Smykalov, P. 2009. School timetabling: Rector. <<http://www.rector.spb.ru/uk/index.html>>
25. Watt. 2009. <<http://www.asap.cs.nott.ac.uk/watt/>>
26. Wren A. 1996. *Scheduling, timetabling and roistering a special relationship?* In Lecture notes in computer science: Vol. 1153. Practice and theory of automated timetabling, pages 46-75. Berlin, Springer.

