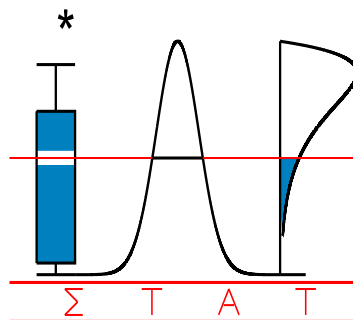


Progress Report 2004  
IAP-network in Statistics  
Contract P5/24

March 28, 2005



# Contents

<b>1</b>	<b>Accomplished research projects</b>	<b>4</b>
1.1	Introduction and overview . . . . .	4
1.1.1	Introduction and overview . . . . .	4
1.1.2	Overview . . . . .	4
1.2	Work package 1: Functional estimation . . . . .	6
1.2.1	Frontier estimation . . . . .	6
1.2.2	Non- and semiparametric regression . . . . .	7
1.2.3	Modelling of heterogeneous regularities . . . . .	10
1.2.4	Inference by means of empirical likelihood techniques . . . . .	10
1.2.5	Functional ANOVA . . . . .	11
1.2.6	Functional estimation for microarray data . . . . .	11
1.2.7	Nonparametric functional estimation by means of wavelets . . . . .	12
1.2.8	Nonparametric density estimation . . . . .	12
1.2.9	Other topics in functional estimation . . . . .	12
1.3	Work package 2: Time series . . . . .	13
1.3.1	Modelling and estimation for non-stationary time series . . . . .	13
1.3.2	Analysis of high-dimensional time series data . . . . .	14
1.3.3	Multivariate time series, spatial data, image analysis, and inverse problems . . . . .	15
1.3.4	Related topics . . . . .	17
1.4	Work package 3: Survival analysis . . . . .	19
1.4.1	Nonparametric estimation with censored data . . . . .	19
1.4.2	Frailty modelling in survival analysis . . . . .	20
1.4.3	Other regression models with censored data . . . . .	20
1.5	Work package 4: Mixed models . . . . .	21
1.5.1	The implementation of multivariate random effects . . . . .	21
1.5.2	The investigation of mixture models as an alternative for approaching the random effects distribution . . . . .	22
1.5.3	Extensions to interval-censored data . . . . .	23
1.5.4	Modeling compliance data . . . . .	25
1.5.5	Mixed models in actuarial sciences . . . . .	25
1.5.6	Functional mixed models . . . . .	26
1.6	Work package 5: Classification and mixture models . . . . .	26
1.6.1	Studying specific types of mixture models . . . . .	26
1.6.2	Investigating methods to decide on the number and type of components . . . . .	28
1.6.3	Classification techniques other than mixtures . . . . .	28
1.6.4	Specific cross-links with other work packages . . . . .	29
1.6.5	Methodological problems . . . . .	29
1.6.6	Functional classification techniques . . . . .	30
1.7	Work package 6: Incompleteness and latent variables . . . . .	30

<b>2</b>	<b>Network activities</b>	<b>32</b>
2.1	Web site . . . . .	32
2.2	Technical reports and published papers . . . . .	32
2.3	Scientific meetings . . . . .	32
2.3.1	Workshops . . . . .	32
2.3.2	Special activities by young researchers . . . . .	33
2.4	Organization of the network : administrative meeting . . . . .	33
2.5	Collaborations, working groups and seminars . . . . .	34
2.5.1	Collaborations . . . . .	34
2.5.2	Working groups . . . . .	34
2.5.3	Seminars . . . . .	35
2.6	Short courses . . . . .	35
2.7	Postdoctoral researchers and return grants . . . . .	36
<b>3</b>	<b>Technical reports and publications</b>	<b>37</b>
3.1	List of publications per research unit/partner . . . . .	37
3.1.1	Université catholique de Louvain, UCL partner . . . . .	37
3.1.2	Katholieke Universiteit Leuven, KUL-1 partner . . . . .	42
3.1.3	Katholieke Universiteit Leuven, KUL-2 partner . . . . .	45
3.1.4	Limburgs Universitair Centrum, LUC partner . . . . .	47
3.1.5	Université Libre de Bruxelles, ULB partner . . . . .	51
3.1.6	Aachen Technical University, RWTH partner . . . . .	55
3.1.7	Université Joseph Fourier, UJF-LMC-IMAG partner . . . . .	55
3.2	List of Joint publications . . . . .	57

# 1 Accomplished research projects

## 1.1 Introduction and overview

### 1.1.1 Introduction and overview

The research project has been built up around six work packages. Table 1 below gives the *main* contributors to each work package and indicates per package the partner that is coordinating the work.

Work package	Contributing partners
WP1: Functional estimation	UCL*, ULB, UJF
WP2: Time series	ULB*, UCL
WP3: Survival analysis	LUC*, UCL
WP4: Mixed models	KUL-2*, KUL-1, LUC
WP5: Classification and mixture models	KUL-1*, KUL-2, RWTH
WP6: Incompleteness and latent variables	LUC*, KUL-1, KUL-2

Table 1: *Main contributors per work package, and coordinating partner per work package (indicated with a \*).*

In the subsections below we describe the progress that has been made in the various work packages. Within each work package we report on the progress that has been achieved on the various *primary objectives* mentioned in the research proposal. For each of the work packages we also indicate **interactions** with research results in other packages: this is done by referring to the other WP as **WP**. The references mentioned in the text can be found in Section 3, which contains a complete list of all publications under the IAP-statistics network.

### 1.1.2 Overview

The overall achievements of the research project can be summarized as follows :

- The work on functional estimation has mainly focused on the following issues. In the context of frontier estimation, further progress has been made on the asymptotic and finite sample properties of the so-called DEA and FDH estimators. The technique of empirical likelihood has been studied in a variety of problems, like for testing in the two-sample problem, and for inference in the presence of a nonparametric nuisance parameter. Researchers of several teams have obtained results in the context of functional ANOVA methods, goodness-of-fit tests for the form of a regression function (both for independent and dependent data), inference for curves under shape restrictions (like monotonicity of a regression or hazard function) and density/regression estimation in random fields.
- In the context of time series, estimation methods for locally stationary processes and for the generalized dynamic factor model had been the main steps in 2002-2003. The year 2004 has been devoted to the crucial forecasting aspects of these two methodologies. A remarkable pointwise adaptive estimator for the spectral density of locally stationary wavelet processes

has also been obtained. The (quasi) likelihood approaches to time-varying models and to dynamic factor models also progressed a lot in 2004. In the spatial domain, two significant contributions on nonparametric estimation have been published, as well as two other ones on inverse problems and inverse imaging. Important results on semiparametric inference techniques based on multivariate ranks and signs to elliptical time series models and shape matrix problems also have been obtained, which resulted in several publications.

- In nonparametric estimation with censored data, new results have been obtained for estimators of bivariate and marginal distributions under dependent censoring. Also estimation of statistical functionals of the form  $Eh(X, Y)$ , where  $X$  is completely observed and  $Y$  is subject to censoring, has been considered. In modelling the presence of heterogeneity in multivariate survival data, important new results have been established via frailties, accelerated failure time and penalized Poisson regression and comparison of these methods has started.
- Among the most important achievements in the area of mixed models is the development of random effects models in accelerated failure time survival models for various types of interval censored data. Allowing a flexible distribution for the error component in combination with a parametric frailty model is a first step towards a general AFT model. Further, the development of pseudo-likelihood methods for highly dimensional mixed effects models is an important step to the use of mixed models in practice. Furthermore the connection between the (generalized) linear mixed models which are very popular in biostatistics and the IRT models brings together two major applied statistical areas. Also, the classical IRT models have been further developed to incorporate in a general way random effects. Finally, correction for misclassification in a logistic (binary, ordinal and count) random effects model has been investigated. The latter is important in epidemiological research.
- In line with previous research, the results on classification and mixture models relate mainly to modeling transition and change with heterogeneous mixture components, smoothing approaches for random effects based on mixture modeling, and the differentiation between categorical and dimensional structures. The latter has resulted in a psychometric publication in the top journal of psychology in general. Concerning classification methods other than mixture models, we were able to develop a comprehensive taxonomy of simultaneous two-way clustering methods, which has led to new models and algorithms. Another development is the generalization of clustering methods with a ‘convexity-based’ clustering criterion. Finally, we have also established an extensive family of models for binary and ordered-category multiway data, in parallel with an existing family of component models for continuous multiway data, and including extensions that go beyond this existing family.
- The work on incompleteness and latent variables has focused on the dissemination and further development of models for incomplete data, with focus on clinical trials, epidemiological studies, and sociological and psychometric application, often in conjunction with the propagation of proper methodology for longitudinal data, including latent structures, mixed-models, and meta-analytic methods. Also, emphasis has been placed on sensitivity assessment tools for incomplete (longitudinal) data. This has resulted in a number of journal and proceedings

publications, as well as book authorships and editorships. Dissemination has also taken place through short courses.

## 1.2 Work package 1: Functional estimation

### 1.2.1 Frontier estimation

This year again, several papers were published or accepted for publication, others have been revised and new working papers have been produced: they cover different approaches and problems involved in this framework.

Most of the work has been devoted to the nonparametric envelopment estimators (DEA and FDH). By construction, these estimators are very sensitive to outliers or extreme values, so we have recently proposed robust versions of envelopment estimators based on concepts of “partial” frontiers (order- $m$  and  $\alpha$ -quantile frontiers). Beguin and Simar (2004) illustrates the use of order- $m$  frontiers in an analysis of the expenses linked to hospital stays to detect outliers. Explaining efficiencies in productivity analysis is also a quite important issue. Daraio and Simar (2005a) propose, in a probabilistic formulation of the problem, a one-stage nonparametric estimator which allows to introduce environmental (or explanatory) variables in the production process. The paper offers also the robust version of the estimator using order- $m$  frontiers. These ideas have been applied in two fields: performance of mutual funds (Daraio and Simar, 2004) and performance of universities (Bonacorsi, Daraio and Simar, 2004). Daraio and Simar (2005b) extend their previous results to convex production sets. Daouia and Simar (2005) present isotonized (monotone) versions of the robust partial frontiers and prove stronger results of convergence. A full multivariate version of  $\alpha$ -quantile frontiers has been proposed in Daouia and Simar (2004), where the asymptotic properties of the nonparametric estimator are established.

Simar and Zelenyuk (2004) develop their research in comparing the efficiency of groups of firms. Here the comparison is in term of a test of the equality of the distribution of the efficiency scores within the two groups. Valentin Zelenyuk (an IAP postdoc at UCL) proposes applications of methodological results obtained in this IAP project: analysis of the catching-up effect in Henderson and Zelenyuk (2004) and the effect of corporate governance in firm’s efficiency in Ukraine, in Zelenyuk and Zheka (2004). He extends his research on aggregation issues in Färe and Zelenyuk (2004).

Envelopment estimators are also by construction biased, Badin and Simar (2004) suggest an easy way for correcting the bias of FDH estimator avoiding the use of the bootstrap. Seok-Oh Jeong is another IAP postdoc at UCL, he analyzes the statistical properties of DEA estimators in Jeong (2004) and Jeong and Park (2004). Jeong and Simar (2005), have investigated the bootstrap algorithm for the FDH and show that the discontinuous nature of the FDH deteriorates its performance in small samples. They provide a linearized version of the FDH which performs much better in finite samples, and a bias-corrected version is also proposed.

“Deterministic” frontier models (as opposed to “stochastic” frontier models) and their non-parametric estimators (DEA and FDH) do not allow for errors or noise in the data. Previous work in this IAP project investigated semi-parametric stochastic models in the presence of panels of data or nonparametric stochastic models restricted to cases where the noise is small with respect

to the signal. A first tentative for approaching nonparametric stochastic frontier models without such restrictions is provided by Kumbhakar, Park, Simar and Tsionas (2004) who propose a local maximum likelihood approach. Here, locally, a parametric model acts as an anchorage model and the localizing procedure allows to fit very flexible models; the results are very promising.

Parametric approaches are also very popular in the field of frontier estimation. Florens and Simar (2005) present a new method based on a parametric approximation of a nonparametric model, which seems to provide a robust parametric estimator, robust to the stochastic hypothesis on the production process but also robust to extreme or outlying observations.

Mouchart and Vandresse (2004) analyse the market imperfection and the bargaining powers of a set of contracts in the sector of the Belgian freight transport. Their method relies on estimating the support of the joint distribution of the price and a set of attributes of each contract; the estimation is non-parametric and is based on an bidirectional extension of the DEA estimator widely used in productivity analysis. This approach is also used to evaluate the role of each attributes in the bargaining process and to assess some possible market segmentation.

At UJF, Girard, Iouditski and Nazin (2005) propose some new optimal estimators for the Lipschitz frontier of a set of points. They are defined as sufficiently regular kernel estimators, covering all the points and whose associated support is of smallest surface.

### 1.2.2 Non- and semiparametric regression

#### Automatic detection of change-points

Gijbels, Lambert and Qiu (2004) deal with nonparametric estimation of a regression curve where the estimation should preserve possible jumps in the curve. At each point  $x$  at which one wants to estimate the regression function, the method chooses in an adaptive way among three different local linear kernel estimates. The choice among these three estimates is made by looking at differences of the weighted residual mean squares of the three fits. The resulting estimate preserves the jumps well and in addition gives smooth estimates of the continuity parts of the curve.

Gijbels and Goderniaux (2004a,b) propose a fully data-driven procedure for detecting jump points in a regression curve or its derivative. They rely on the two-steps estimation method. Gijbels and Goderniaux (2004a) also discuss how to estimate the number of jump points in a regression curve, using cross-validation. They also provide a fully data-driven algorithm that should be used when dealing with an identification problem, that typically occurs when a jump point is difficult to distinguish from points with a high derivative (in absolute value). The basic methodology can be adopted to the case of jump points in a derivative of the regression curve. This requires some adjustments such as the choice of an appropriate diagnostic function and a cross-validation criterion for derivative estimation.

Of interest is also the testing problem. Gijbels and Goderniaux (2004c) dealt with testing the null hypothesis that a regression is continuous versus the alternative hypothesis that it is a discontinuous function, relying on the two-steps data-driven procedure. In addition they introduce a bootstrap procedure for assessing the distribution of the test statistic under the null hypothesis. This bootstrap procedure has also been used in Gijbels, Hall and Kneip (2004) for constructing confidence bands for discontinuous regression curves. The testing procedure proposed by Gijbels and Goderniaux (2004c) has been compared with other testing procedures available in the liter-

ature. The latter procedures rely on the use of the asymptotic distribution of the involved test statistics, and do not treat in a satisfactory way the choice of the smoothing parameters involved.

### **Goodness-of-fit tests**

Penalized regression spline models afford a simple mixed model representation where variance components control the degree of non-linearity in the smooth function estimates. This is the motivation for Claeskens (2004) to study lack of fit tests based on the restricted maximum likelihood ratio statistic. Test statistics are studied for simple as well as multiple regression models.

In Claeskens and Hjort (2004) classes of goodness-of-fit tests are constructed via estimated likelihood ratios and use log-linear expansions. Such expansions are either coupled with subset selectors like the AIC and the BIC, or use order growing with sample size. Our tests are generalized to testing adequacy of general parametric models, and work also in higher dimensions.

Aerts, Claeskens and Hart (2004) propose and analyze nonparametric tests of the null hypothesis that a regression function belongs to a specified parametric family. The tests are based on BIC approximations to the posterior probability of the null model, and may be carried out in either Bayesian or frequentist fashion. Simulation results and an example involving variable star data illustrate desirable features of the proposed tests.

Van Keilegom, González-Manteiga and Sánchez-Sellero (2004) also construct a test statistic for the hypothesis that a regression function belongs to some parametric family, but their approach is based on a different idea. Their test statistic is based on the distance between the empirical distribution function of the parametric and of the nonparametric residuals. They construct classical Kolmogorov-Smirnov and Cramér-von Mises statistics based on this distance. The same authors are currently working on an empirical likelihood based test for the same null hypothesis. It is hoped that this test will lead to good power in comparison with its competitors, due to the powerful empirical likelihood methodology.

All of the above papers are restricted to the case of independent observations. When the errors form a stationary time-series, Wang and Van Keilegom (2004) propose a new nonparametric method for testing the parametric form of a regression function. The nonparametric test is motivated by recent advancement in the theory of ANOVA with large number of factor levels. The paper forms a bridge between WP1 and **WP2**.

A problem related to the ones above is the problem of testing whether two populations have the same regression curve. In Pardo-Fernández, Van Keilegom and González-Manteiga (2004), this hypothesis is tested by using an approach similar in spirit to the one developed by Van Keilegom, González-Manteiga and Sánchez-Sellero (2004) described above, i.e. their approach is also based on the comparison of the empirical distribution of the residuals under the null hypothesis, and the residuals obtained without using the null hypothesis. Pardo-Fernández and Van Keilegom are currently working on an extension of this test to the case of censored data; this paper will be strongly related to **WP3**.

### **Inference for curves under shape restrictions**

Antoniadis, Bigot and Gijbels (2005) (J. Bigot was a IAP postdoc at UCL in 2004 and got his PhD at UJF) focus on nonparametric estimation of a constrained regression function using penalized wavelet regression techniques. This results into a convex optimization problem under linear con-



straints. The estimator is easily obtained via the dual formulation of the optimization problem. In particular, a penalized wavelet monotone regression estimator is investigated. Rate of convergence and finite sample performances of this estimator are also studied.

In many applications it is reasonable to assume that an unknown function satisfies certain qualitative constraints, such as monotonicity. A review on procedures for testing for monotonicity of a regression function is provided in Gijbels (2004). Gijbels and Heckman (2004) deal with the problem of testing whether the hazard function can indeed be assumed to be increasing (or decreasing). They also illustrate how to apply their methodology to type II censored data. The paper is therefore also linked to **WP3**.

A related problem is studied in Hall and Van Keilegom (2004). They construct a test for the null hypothesis that a hazard rate is monotone nondecreasing, versus the alternative that it is not. Both the test statistic and the means of calibrating it are new. Unlike previous approaches, neither is based on the assumption that the null distribution is exponential.

### **Estimation of the regression curve**

Amato, Antoniadis and Pensky (2004) consider regression problems with univariate design points. The design points are irregular and no assumptions on their distribution are imposed. The regression function is retrieved by a wavelet based reproducing kernel Hilbert space (RKHS) technique with the penalty equal to the sum of blockwise RKHS norms. Under additional assumptions on design points the method achieves asymptotic optimality in a wide range of Besov spaces.

Antoniadis, Gijbels and Nikolova (2005) study a general methodology for nonparametric regression via regularization. The approach chosen consists of penalized likelihood regression for generalized linear models with nonquadratic penalties.

The related problems of estimating density and regression functions in the spatial context of random fields have been considered in Hallin, Lu, and Tran (2004a and b); see **WP2**, Section 1.3.3 for details.

### **Dimension reduction techniques**

Amato, Antoniadis and De Feis (2004) develop some new two-dimensional reduction regression methods to predict a scalar response from a discretized sample path of a continuous time covariate process. The methods take into account the functional nature of the predictor and are both based on appropriate wavelet decompositions. Using such decompositions, they derive prediction methods that are similar to minimum average variance estimation (MAVE) or functional sliced inverse regression (FSIR).

Single-index models offer a flexible semiparametric regression framework for high-dimensional predictors. Bayesian methods have never been proposed for such models. Antoniadis, Grégoire and McKeague (2004) develop a Bayesian approach incorporating some frequentist methods: B-splines approximate the link function, the prior on the index vector is Fisher von Mises, and regularization with generalized cross validation is adopted to avoid over-fitting the link function. A random walk Metropolis algorithm is used to sample from the posterior.

An extension of PLS and dimension reduction in logit models is proposed in Fort and Lambert-Lacroix (2004a). The extension works also when the number of covariates is far larger than the number of observations.

Fort, Lambert-Lacroix and Peyre (2004) propose nonparametric single-index models in generalized linear models as potential reduction method for supervised classification of microarray data.

G. Geenens in his PhD (advisor: L. Simar) analyzes, in collaboration with M. Delecroix (ENSAI, Rennes) the possibility of modelling the probabilities in a multinomial process by single index models (SIM). The idea is to let these probabilities be related to some explanatory variables in a semiparametric way. All the techniques of estimation for SIM have been investigated in this particular setup. Then these methods will be used for doing inference in contingency tables (test of independence, of partial independence, . . .). Geenens and Delecroix (2005) developed a survey on SIM models, while Geenens, Simar and Delecroix are currently investigating the properties of a new estimator based on a maximum rank correlation criterion.

### **Singularity estimation in regression**

Bigot (2005) deals with the problem of determining the singularities of a curve when it is observed with noise. A nonparametric approach, based on appropriate thresholding of the empirical continuous wavelet coefficients, is proposed to estimate the wavelet maxima of a noisy signal. A new tool, “the structural intensity”, is also introduced to represent the locations of the singularities of an unknown signal via a density function. This approach is shown to be an effective technique for detecting the significant singularities of a noisy signal and for removing spurious estimates.

### **1.2.3 Modelling of heterogeneous regularities**

Delouille, Simoens and von Sachs (2004) have presented a new methodology to construct smooth wavelets which adapt automatically to the stochastic design of a non-parametric curve estimation problem, and which circumvent the usual restrictions of classical wavelets (dyadic sample size, boundary treatment, non-equispaced design). Delouille and von Sachs (2004) furnish some theoretical results on the optimality of this new estimator which parallel classical wavelet thresholding for equispaced data. Moreover they develop these in the more difficult time series context of a non-linear autoregressive design and show how their original algorithm can be adapted to various time series situations (including ARCH-type models). This research is linked with research under **WP2**. For the two-dimensional situation, Delouille, Jansen and von Sachs (2004) present an approach to treat fully irregularly spaced data in two dimensions. Smooth bivariate estimators are constructed based on lifting on triangulation schemes using Lagrange interpolating polynomials and a Bayesian approach for wavelet thresholding.

### **1.2.4 Inference by means of empirical likelihood techniques**

Empirical likelihood methods offer an alternative to the classical approaches based on asymptotic normality or bootstrap approximation. In Hjort, McKeague and Van Keilegom (2004) (conditionally accepted for *The Annals of Statistics*) the scope of the empirical likelihood methodology is extended in three directions: (1) to allow for plug-in estimates of nuisance parameters in estimating equations, (2) slower than  $\sqrt{n}$ -rates of convergence, and (3) settings in which there are a relatively large number of estimating equations compared to the sample size. Calibrating empirical likelihood confidence regions with plug-in is sometimes intractable due to the complexity of the

asymptotics, so they introduce a bootstrap approximation that can be used in such situations. A range of examples from survival analysis (see **WP3**) and nonparametric statistics are provided to illustrate the main results. Although not considered in the paper, the results can be applied to a much wider class of applications, including time series models (see **WP2**), and settings with incomplete data (see **WP6**).

Cao and Van Keilegom (2004) study the problem of testing whether two populations have the same law, by comparing kernel estimators of the two density functions. The proposed test statistic is based on a local empirical likelihood approach. They obtain the asymptotic distribution of the test statistic and propose a bootstrap approximation to calibrate the test. An extension to functional data is in preparation, in collaboration with Peter Hall.

Finally, a paper on U-quantile estimation in the presence of auxiliary information by means of empirical likelihood techniques is in preparation by Van Keilegom and Veraverbeke (2005).

### 1.2.5 Functional ANOVA

In Bugli and Lambert (2004) (UCL) it is shown how one can use P-splines in a functional ANOVA framework to decompose event related potentials extracted from electroencephalograms into physiologically meaningful components from which the effect of a drug on the brain can be interpreted and quantified.

At UJF, Abramovich, Antoniadis, Spatinas and Vidakovic (2004) consider the testing problem in a fixed-effects functional analysis of variance model. We test the null hypotheses that the functional main effects and the functional interactions are zeros against the composite nonparametric alternative hypotheses that they are separated away from zero in L2-norm and also possess some smoothness properties. The corresponding tests are based on the empirical wavelet coefficients of the data.

### 1.2.6 Functional estimation for microarray data

Mercier, Berthault, Mary, Peyre, Antoniadis, Comet, Cornuejol, Froidevaux and Dutreix (2004) combine two independent analysis methods (ANOVA and RELIEF) to compare gene expression patterns in *Saccharomyces cerevisiae* growing in the absence and continuous presence of varying low doses of radiation. Global distribution analysis highlights the importance of mitochondrial membrane functions in the response. It is demonstrated that microarrays detect cellular changes induced by irradiation at doses that are 1000-fold lower than the minimal dose associated with mutagenic effects.

Fort and Lambert-Lacroix (2004b) propose a new method combining Partial Least Squares and Ridge penalized logistic regression for adaptive classification of microarray data. Their interest is outlined in some cases, and explain theoretically the poor behavior of some currently used classifiers. The proposed procedures are compared with these other classifiers and the predictive performance of the resulting classification rule is illustrated on two well known data sets: the Leukemia data set and the Colon data set.

### 1.2.7 Nonparametric functional estimation by means of wavelets

Wavelet-based denoising techniques are well suited to estimate spatially inhomogeneous signals. Waveshrink (Donoho and Johnstone) assumes independent Gaussian errors and equispaced sampling of the signal. In Sardy, Antoniadis and Tseng (2004) a unifying L1-penalized likelihood approach to regularize the maximum likelihood estimation by adding an L1 penalty of the wavelet coefficients is developed. The approach works for all types of wavelets and for a range of noise distributions.

Wavelet thresholding is a natural and efficient approach for removing noise from data in nonparametric function estimation. In order to achieve spatial adaptation and benefit from the sparse representation of most signals in the wavelet domain, it is crucial to choose the thresholds correctly. An important property of the universal thresholding is that the reconstruction is noise-free, i.e. when the true signal is constant, then, with high probability, the estimated function is also constant and equal to the empirical mean of the data. Motivated by this noise-free reconstruction property, Antoniadis and Fryzlewicz (2004) investigate a parametric thresholding procedure which takes advantage of the increasing sparsity of the wavelet coefficients across scales.

Antoniadis and Bigot (2004) focus on nonparametric estimators in inverse problems for Poisson processes involving the use of wavelet decompositions. Their approach combines Galerkin inversion methods with the use of log-intensity functions approximated by wavelet series. This method inherits the well known theoretical advantages of wavelet-vaguelette decompositions for inverse problems together with the remarkably simple closed form expressions of Galerkin approximations.

### 1.2.8 Nonparametric density estimation

Bouezmarni and Rolin (2004) (T. Bouezmarni is currently IAP-postdoc at KUL-1, but obtained his PhD at UCL) give the exact asymptotic behaviour of the expected average absolute error of a beta kernel density estimator. They also prove the uniform weak consistency of this estimator for the class of continuous densities.

Delaigle and Gijbels (2004a,b) deal with the problem of how to estimate nonparametrically the density of a random variable when the measurements on this variable contain errors. They propose practical bandwidth selection procedures. A finite sample comparison between the discussed procedures, a bootstrap procedure, a plug-in procedure and a cross-validation procedure, has been carried out. Theoretical properties (including consistency) of the bootstrap and plug-in procedures have been established.

### 1.2.9 Other topics in functional estimation

#### Backfitting/profiling in semiparametric models

Van Keilegom and Carroll (2005) study the backfitting and profile methods for general criterion functions that depend on a parameter of interest and a nuisance function. They show that when different amounts of smoothing are employed for each method to estimate the nuisance function, the two estimation procedures produce asymptotically the same estimator of the parameter of interest. The results are applied to a partial linear median regression model and a change point model. Although not considered in the paper, the results could also be applied to models from

other domains of applications, including survival analysis (**WP3**).

### **Extreme value methods**

Gardes and Girard (2005) present and study a new estimator of the extreme value index of a distribution adapted to any domain of attraction. Its construction is similar to the one of Pickand's estimator.

### **Aggregation**

Hoeffelman et al. (2004) propose a methodology to derive a statistical distribution to represent the magnetic field observed over time below HV overhead power lines in Belgium. It is based on a spatial aggregation of local distributions estimated from magnetic field measured below a representative sample of spans of the Belgium power lines network. Two possible parametric or empirical methods are used to estimate local distribution and are combined with the known load distribution over the year to derive the space-time aggregated distribution. Summary parameters like mean, percentiles and risks levels are derived from the aggregated distribution as well as corresponding bootstrap confidence intervals.

### **Time intensity curves**

Guyot et al. (2005) are presenting the results of a study of beer astringency based on a sensory analysis through time intensity curves. A methodology is set up to analyse the designed experiment combining a dimension reduction of intensity curves in some queue indexes and the analysis of resulting responses using polynomial mixed models to take into account important judge effects.

### **Copula modeling**

In Vandenhende and Lambert (2004) a new family of Archimedean copulas is defined using a continuous piecewise log-linear combination of existing Archimedean generators. We provide an efficient least-squares method to estimate the involved coefficients. We show that these coefficients can be interpreted as local dependence measures. A smooth estimate is obtained with a penalized approach.

Denuit, Purcaru and Van Keilegom (2004) consider an application of bivariate Archimedean copula modeling in non-life insurance. The data in their application are subject to censoring (an observation is censored when the amount of the claim exceeds the policy limit). The paper has a natural connection with **WP3**.

### **Poisson regression**

Brouhns, Denuit and Van Keilegom (2004) describe the Poisson log-bilinear model for mortality projection. They derive confidence intervals for expected remaining lifetimes with the help of the bootstrap. Belgian mortality statistics are investigated using these techniques.

## **1.3 Work package 2: Time series**

### **1.3.1 Modelling and estimation for non-stationary time series**

The study of this topic has been pursued, mainly by R. von Sachs and S. Van Bellegem at UCL (mainly, using a local stationarity approach), and by G. Méléard and his collaborators at ULB (in a time-varying coefficient perspective). M. Hallin and S. Lotfi have investigated some properties

of multivariate time series models with periodic coefficients.

### **Locally stationary processes**

In Van Bellegem and von Sachs (2004a), emphasis is focussed on a simple and meaningful model for variance nonstationarity. This model satisfactorily explains the nonstationary behavior of several economic data sets, among which are the U.S. stock returns and exchange rates. The question of how to forecast these processes is addressed and evaluated on the data sets. In Van Bellegem and von Sachs (2004b), the same authors present an algorithm for the pointwise adaptive estimation of the spectral density of locally stationary wavelet processes. As an application they derive a new test of covariance stationarity and a test of the local significance of the coefficients of the possibly sparse wavelet representation of the underlying model of locally stationary wavelet processes.

Ombao, von Sachs and Guo (2004) address the problem of complexity reduction for high-dimensional second-order non-stationary data such as EEG with transient phenomena. An automatic, fast and fully nonparametric algorithm, based on principal component analysis in the frequency domain, is proposed to adapt a piecewise stationary frequency domain model to the data and to estimate the spectral density matrix, phase and coherency in this model.

### **Time-varying coefficient models**

The paper by Azrak and M elard (2004) on quasi-maximum likelihood estimation of the parameters of ARMA models with time-dependent coefficients is now to appear in *Statistical Inference for Stochastic Processes*. Two other papers are in progress: a generalization to multidimensional processes and a paper focussing on pure AR models, including a comparison with the Dahlhaus approach.

S. Lotfi has defended her thesis (supervisor: M. Hallin), with a dissertation entitled “Efficient tests for the periodic structure of some time series models”. Together with M. Hallin, she has published a paper (Hallin and Lotfi, 2004) on the optimal detection of periodicities in the coefficients of  $m$ -variate  $d$ -periodic VAR( $p$ ) models. VAR models with periodic coefficients are an attractive alternative to the more traditional differencing approach to seasonality.

## **1.3.2 Analysis of high-dimensional time series data**

### **The dynamic factor model as a tool for forecasting macroeconomic time series with high dimensional data**

Investigation of the properties, applications, and extensions of the generalized dynamic factor method has been pursued. Two important papers in that area have appeared. The first one (Forni, Hallin, Lippi, and Reichlin (2004)) deals with the delicate problem of rates of consistency (as both the number  $n$  of series, and their length  $T$  tend to infinity) of the spectral estimators derived in the earlier papers by the same authors. The second one (Forni, Hallin, Lippi, and Reichlin (2005)) considers the forecasting problem for very high dimensional time series. This paper proposes a new forecasting method taking advantage, via the dynamic factor model, of the information on the dynamic covariance structure of the whole panel. We first obtain an estimation for the covariance matrices of common and idiosyncratic components. The generalized eigenvectors of this couple of matrices are then used to derive a consistent estimate of the optimal forecast. This two-step approach solves the end-of-sample problems caused by two-sided filtering, while retaining

the advantages of an estimator based on dynamic information. The relative merits of this method and the one proposed earlier by Stock and Watson are carefully discussed.

The project described in last year report, conducted with the European Central Bank and the Board of Governor at the Federal Reserve are now internal publications, not yet available for circulation. This work aims at identifying the roles of different blocks of data releases for the forecast of output and inflation. It develops a method, based on the use of the Kalman filter, to handle non-synchronized intra-month releases. This work is conducted in collaboration by Lucrezia Reichlin and Domenico Giannone, and will soon be available in the form of working papers.

D'Agostino and Giannone (2005) coauthored a paper examining parameterization of the factor model in an out-of-sample forecasting exercise.

### **Structural forecasting using large dimensional datasets**

This work aims at identifying shocks computed from the factor model forecast using restrictions based on economic theory. The paper by Giannone, Reichlin and Sala (2005a) is now completed. A similar topic is treated in Forni, Giannone, Lippi and Reichlin (2004). This is a revision of Forni, Lippi and Reichlin (2003). A new part showing rates of convergence for the estimates of the shocks and their coefficients is added. Giannone, Reichlin and Sala (2005b) also analyze factor models as tools for estimating general equilibrium macroeconomic models.

### **Quasi maximum likelihood estimation for structural factor models**

The work that was started last year by L. Reichlin with C. Doz and D. Giannone on linking principal components and maximum likelihood estimation in factor models for large cross-sections, is in its final stage, and soon will concretize into a manuscript.

### **Business cycle research using a variety of techniques**

Reichlin has edited a book on the Euro area business cycle (Reichlin, 2005) and co-authors a paper in the volume (Giannone and Reichlin, 2005). Several projects are in progress.

## **1.3.3 Multivariate time series, spatial data, image analysis, and inverse problems**

### **Semi-parametric rank-based approach to the analysis of multivariate time series**

A series of papers by M. Hallin and D. Paindaveine, dealing with rank-based inference for multivariate time series models with elliptical innovation densities, has appeared (2004a,b, 2005a,b,c), where optimal (in the Le Cam sense) distribution-free tests are based on multivariate concepts of signs and ranks. One of those concepts of multivariate ranks, based on hyperplane counts, has been introduced in Oja and Paindaveine (2005). This concept is an analogue, for distances between the observations, of the so-called Randles interdirections associated with the observed angles. Based on these concepts, these papers mainly develop optimal rank-based procedures for testing affine-invariant linear hypotheses on the parameters of a multivariate general linear model with elliptical VARMA errors. A class of optimal procedures is proposed, that generalize the univariate signed rank procedures proposed in the literature by Hallin and Puri, and are locally asymptotically most stringent under correctly specified radial densities. Their AREs with respect to Gaussian procedures are shown to be convex linear combinations of the AREs obtained in 2002 by Hallin and Paindaveine for the pure location and purely serial models, respectively. The result-

ing test statistics are provided under closed form for several important particular cases, including generalized Durbin-Watson tests, VARMA order identification tests, etc. In Hallin and Paindaveine (2005b) these general results are put at work in two examples of practical relevance: (i) the multivariate Durbin-Watson problem (testing against autocorrelated noise in a linear model context), and (ii) the problem of testing the order of a vector autoregressive model (testing  $\text{VAR}(p_0)$  against  $\text{VAR}(p_0 + 1)$  dependence). These two testing procedures are the building blocks of classical autoregressive order-identification methods. Based either on pseudo-Mahalanobis (Tyler) or on hyperplane-based (Oja and Paindaveine (2005)) signs and ranks, three classes of test statistics are considered for each problem: (a) statistics of the sign test type, (b) Spearman statistics, and (c) van der Waerden (normal score) ones. Simulations confirm theoretical results about the power of the proposed rank-based methods, and establish their good robustness properties.

### **Image analysis and inverse problems**

During the year 2004, Christine De Mol has pursued her work on the use of sparsity-enforcing penalties for the regularization of ill-posed linear inverse problems as well as for nonparametric regression. A typical penalty of this type is the  $\ell_1$  norm of the sequence of coefficients of the expansion of the solution on a given arbitrary orthonormal basis. In Daubechies, Defrise and De Mol (2004), it is shown that such penalty regularizes the problem and an iterative algorithm is devised to compute the corresponding solutions. The results hold more generally when the  $\ell_1$  norm is replaced by a weighted  $\ell_p$  norm with  $1 < p < 2$ , the limit case  $p = 2$  corresponding to the usual quadratic regularization applied for example in ridge regression. Defrise and De Mol (2004) have generalized this framework to cover the case of mixed penalties. Using the technique of surrogate functionals, new iterative algorithms have been derived for solving linear inverse problems of this type which are encountered in several applications. As a special case, a Huber penalty or prior can be implemented, i.e. a quadratic penalty on the small coefficients and a  $\ell_1$  penalty on the larger ones. On the other hand, in connection with the Ph. D. project of M. Banbura, C. De Mol has been working on the analysis and forecasting of nonstationary time-series, with emphasis on economic data. In the light of inverse problem theory, several methods proposed in the literature can be reformulated in a unified framework. Several methods have been tested and compared numerically, with as a benchmark the inflation series in the univariate case and several macroeconomic panels in the multivariate case. Wavelet-based analysis and forecasting methods have also been devised and are under current investigation.

### **Spatial data and the analysis of random fields**

M. Hallin has pursued his collaboration on random fields with Z. Lu and L.T. Tran. This collaboration has produced two publications in 2004. In the first one, Hallin, Lu and Tran (2004a) investigate kernel density estimators for spatial processes with linear or nonlinear structures. Sufficient conditions for kernel estimators to converge in  $L_1$  are obtained under extremely general, verifiable conditions. The results hold for mixing as well as for non-mixing processes. Potential applications include testing for spatial interaction, the spatial analysis of causality structures, the definition of leading/lagging sites, the construction of clusters of comoving sites, etc. In the second paper (Hallin, Lu and Tran (2004b)), a local linear kernel estimator of the regression function of a stationary spatial process observed over a rectangular domain is proposed and investigated. Under



mild regularity assumptions, asymptotic normality of the estimators of the regression function and its derivatives is established. Appropriate choices of the bandwidths are proposed. The spatial process is assumed to satisfy some very general mixing conditions, generalizing classical time-series strong mixing concepts. The size of the rectangular domain is allowed to tend to infinity at different rates.

A spatial regression quantile approach to the same problem is under progress. This approach yields a much richer information, with an analysis of the spatial conditional quantile structure. To the best of our knowledge, this is the first time regression quantiles are considered in a spatial context.

### **Independence between multivariate time series**

The problem of testing non-correlation and non-causality (in the Granger sense) between two multivariate series is investigated in Hallin and Saidi (2004a,b). The paper by Hallin and Saidi (2004a) develops an approach to the problem of testing non-correlation (at all leads and lags) between two univariate time series. Their test generalizes the test earlier proposed by Koch and Yang to the multivariate case. A Monte Carlo study is conducted, which indicates that their approach performs better than a test proposed by El Himdi and Roy for a wide range of models. Both procedures are applied to the problem of testing the absence of correlation between Canadian and U.S. economic indicators, and to a brief study of causality between money and income in Canada.

The second paper (Hallin and Saidi (2004b)) considers the same problem from the point of view of local and asymptotic optimality. Assuming that the global process  $\{\mathbf{X}_t, t \in Z\} := \{((\mathbf{X}_t^{(1)})^T, (\mathbf{X}_t^{(2)})^T)^T, t \in Z\}$  admits a joint vector autoregressive (VAR) representation, they construct locally and asymptotically optimal pseudo-Gaussian tests for the null hypothesis of non-correlation between  $\{\mathbf{X}_t^{(1)}\}$  and  $\{\mathbf{X}_t^{(2)}\}$ , and compare their local asymptotic powers with those of the various tests (Haugh-El Himdi-Roy, and Koch-Yang-Hallin-Saidi) proposed in the literature.

Paindaveine (2004) establishes Chernoff-Savage and Hodges-Lehmann results for the related problem of testing independence between two multivariate iid sequences of observations.

### **1.3.4 Related topics**

#### **Computational issues in time series estimation**

Together with A. Klein, J. Niemczyk, T. Zahaf and P. Spreij, G. Mélard has pursued his investigations on the information matrix of time series models. A paper by Klein and Mélard (2004) has been published on the computation of the asymptotic Fisher information of univariate SISO (single input, single output) models where the polynomials appear as factors of regular and seasonal polynomials. A note by Klein, Mélard and Niemczyk (2004) on the exact information matrix of a multivariate Gaussian process has been submitted, building on a previous paper by Klein, Mélard and Zahaf. A paper by Klein, Mélard and Spreij (2005) on the asymptotic information matrix of a VARMA model has now been accepted for publication. This latter paper establishes a new algebraic characterization of VARMA models by means of a tensor Sylvester matrix, and shows how the Fisher information matrix evaluated numerically, contrary to the tensor Sylvester matrix, may fail to reveal common eigenvalues in the AR and MA matrix polynomials.

The paper by Mélard, Roy and Saidi (2005) on the computational estimation of the parameters

of VARMA models either in structured form (scalar component model or echelon form) or with a unit root (the case of a partially non stationary model) is tentatively accepted, subject to minor revision.

A. Ouakasse has defended his thesis (supervisor: G. Mélard) on a new method for recursive estimation of ARMA and SISO models. Two papers in collaboration with G. Mélard are being completed: one on ARMA models with a new application, and one on the SISO model. A third one, joint with T. Zahaf, is planned.

Finally, Niemczyk (2004) considers the derivatives of the autocovariances of a VARMA process (with respect to the VARMA coefficients), and studies two applications.

### **Statistical inference for shape matrices**

M. Hallin, D. Paindaveine and H. Oja have devoted two papers (Hallin and Paindaveine 2005c; Hallin, Oja and Paindaveine 2004) to optimal rank-based inference methods (tests and estimation) for the shape matrix of an elliptically contoured multivariate density. Contrary to the everyday practice Gaussian procedures, which rely on empirical covariance matrices and require finite moments of order four, these methods remain valid ( $\alpha$  level tests and root- $n$  consistent estimators) without any finite moments (even of arbitrarily small order  $\delta > 0$ ). Paindaveine (2004) moreover shows that the celebrated Chernoff-Savage result establishing the uniform superiority (in the Pitman sense) of normal-score rank-based methods for location over their pseudo-Gaussian competitors also holds for shape matrices.

The rank tests and the R-estimators derived in the two other papers thus uniformly dominate everyday practice, which moreover requires finite fourth order moments.

### **Rank-based and distribution-free inference for time series models**

M. Hallin, in collaboration with C. Vermandele and B. Werker, also devotes two papers to rank-and-sign techniques for serial and non-serial median-restricted models (Hallin, Vermandele and Werker (2003, 2004)). These techniques improve on more standard LAD techniques by considering, along with the signs, the ranks of the observations. Both hypothesis testing and estimation problems are considered. Such techniques rely on a new class of statistics, involving the signs and the ranks (not the *signed ranks*, which require symmetry), which are maximal invariant and distribution-free under median-centered white noise. It is shown that semiparametric efficiency always can be reached by means of these new statistics.

Hallin, Jurečková, and Koul (2005) introduce regression and autoregression rank score statistics of the serial type which, contrary to those existing in the literature, are completely measurable with respect to the (auto)regression rank scores (the existing ones are mixing (auto)regression rank scores with past observations). They establish an asymptotic representation result, and the asymptotic normality of these new statistics, and show how they can be used as a tool for inference in (univariate) time series. The remarkable feature of these genuinely (auto)regression rank score-based statistics is that, unlike their aligned rank counterparts, they are totally insensitive to any shift effect related with the estimation of underlying nuisance parameters.

Dufour, Fahrat and Hallin (2005) establish exact bounds for the tail areas of distributions of autocorrelation coefficients under unspecified innovation density. These bounds allow for performing finite-sample conservative tests irrespective of the underlying symmetric innovation densities.

## Other topics in time series

A. Ayadi and G. Mélard (2004) studied the (very) small sample behavior of autocorrelations of a Gaussian white noise process, with exact results and simulation results in the case of some non-Gaussian white noise processes.

Prediction from panel data with a high proportion of missing data raises substantial statistical issues. Mouchart and Rombouts (2005) propose an efficient approach for the case of nowcasting, i.e. forecasting present values based on recent past data. A progressive specification strategy is elaborated and illustrated on R&D data for a panel of countries from the European Community.

H. Njimi under the supervision of G. Mélard is working on the improvement of the automatic ARMA modeling procedure of Mélard and Pasteels. In order to apply the methodology to the shortest series of the M3-Competition, he has considered using Ayadi and Mélard's work on autocorrelations. The paper by Azrak, Mélard and Njimi (2004) where this strategy is applied to a problematic data set with many missing values and outliers has been published. Meanwhile, he is progressing on a comparison with the Gómez and Maravall TRAMO/SEATS method.

## 1.4 Work package 3: Survival analysis

### 1.4.1 Nonparametric estimation with censored data

At LUC, Veraverbeke (2004) studied nonparametric estimation of two important functionals of the conditional residual lifetime beyond some fixed or random timepoint: the mean and any quantile. The observations are subject to random censoring and covariate information is taken into account. Paul Janssen and Noël Veraverbeke, in collaboration with Jan Swanepoel (Potchefstroom, South Africa) studied variable bandwidth kernel estimators for distribution functions. Their new estimator reduces the bias considerably and keeps the variance unchanged with respect to the usual kernel distribution function estimator. See Janssen, Swanepoel and Veraverbeke (2004). Li Chun Wang proposed a criterion for choosing between two loss functions in a Bayesian analysis. See Wang (2004a). He also studied Bayes and empirical Bayes tests for the life parameter and studied asymptotically optimal rates of convergence. See Wang (2004b). Together with Veraverbeke he studied Bayes prediction in the exponential distribution under random censorship and where the prior distribution is of unknown form. See Wang and Veraverbeke (2004a). In Wang and Veraverbeke (2004b) they define an adjusted empirical log-likelihood for the mean response in the situation where responses are missing at random. This last paper has natural links with **WP1** (subsection 1.2.4 on empirical likelihood) and **WP6**.

At UCL, Van Keilegom (2004) proposes estimators of the bivariate and marginal distributions of two variables subject to censoring. The estimators do not require the common assumption of independence between the vector of survival and censoring times, but allow for a certain type of dependent censoring. As an application she discusses the estimation of the regression coefficients in a polynomial regression model, when both the response and the covariate are subject to censoring.

### 1.4.2 Frailty modelling in survival analysis

In 2004 the work on modelling multivariate survival data focused on the following issues. In collaboration with the European Organisation for Research and Treatment of Cancer (EORTC, Brussels) frailty models have been used to study the validity of prognostic indices and to understand heterogeneity in outcome between participating centres in multicentre studies. To support these studies, Legrand, Ducrocq, Janssen, Sylvester and Duchateau (2004) and Legrand, Duchateau, Sylvester, Janssen, van der Hage, van de Velde and Therasse (2004) extended the Bayesian approach in Survival Kit to include a random treatment-centre interaction effect. Duchateau, Opsomer, Dewulf and Janssen (2005) continued the work on the use of splines to arrive at more flexible frailty models. A further important topic is to find appropriate ways to explain the presence of heterogeneity in multivariate survival data in terms of quantities that are relevant and understandable for the users, e.g., how does heterogeneity influence the median survival time or the probability that the disease free survival is at least five years (see Duchateau and Janssen (2004)). Massonnet, Burzykowski and Janssen (2004) (LUC) continued the work on resampling plans for frailty models and they are exploring a new method to estimate the heterogeneity parameters in frailty models via model transformations. The idea is that after the model has been transformed they can rely on the well established mixed model methodology to estimate the heterogeneity parameters. To model the presence of heterogeneity in multivariate survival data different approaches have been proposed in the statistical literature. Frailty models are an approach that adds a random effect to the Cox regression model. Other approaches include accelerated failure time models for correlated survival models (Lesaffre and Komarek (KUL-2)) and the penalized Poisson regression approach (Lambert (UCL)). An important question is to study whether these three different methodological tools lead to the same conclusions when applied to the same data set. To answer this question they are analyzing the data of an EORTC multicentre trial; this will lead to a comparison of the different models that are used to do statistical inference on heterogeneity. Duchateau and Janssen recently signed a contract with Springer to write a book on frailty models. This work started October 2004 and should be finished in the summer of 2006.

### 1.4.3 Other regression models with censored data

Accelerated failure time models with a shared random component are described in Lambert, Collett, Kimber and Johnson (2004) and they are used to evaluate the effect of explanatory factors and different transplant centres on survival times following kidney transplantation. Different combinations of the distribution of the random effects and baseline hazard function are considered and the fit of such models is critically assessed. A mixture model that combines short-term and long-term components of a hazard function is then developed.

In Lambert and Eilers (2004) it is shown how one can use Poisson log-linear models in combination with Bayesian P-splines to set up flexible models for the hazard rate with time-varying regression coefficients. Bayesian inference tools based on the Metropolis-adjusted Langevin algorithm are proposed.

For arbitrary functions  $\varphi$ , Sánchez-Sellero, González-Manteiga and Van Keilegom (2004) consider the problem of estimating expectations of the form  $E[\varphi(X, Y)]$  (where  $X$  is completely ob-

served and  $Y$  is subject to random censoring), which appear in many statistical problems. Applications to goodness-of-fit testing in regression (see also **WP1**, Section 1.2.2) and to inference for the regression depth, are considered in more detail.

## 1.5 Work package 4: Mixed models

The research topics treated in this work package are subdivided into three related themes showing also a high relationship with other work packages.

### 1.5.1 The implementation of multivariate random effects

The following topics were treated:

#### **Flexible distributions for the random effects part of a linear mixed model**

Members of the KUL-2 team (Ghidey and Lesaffre) have developed, in collaboration with Paul Eilers a mixed model with a smooth random effects distribution. The model assumes a flexible random effects distribution that can be well approximated by a smooth function of B-splines or of Gaussian densities. Penalized likelihood maximisation delivers the estimated fixed effects and the smoothing coefficients of the random effects distribution. This work has been published in *Biometrics*, see Ghidey, Lesaffre and Eilers (2004). An extensive simulation study has been conducted (and almost finalized) to compare the performance of the approach to its competitors in the case when other assumptions of the linear mixed model (on the error part and with regard to the process of missing data) are not fulfilled. The simulation results indicate superiority of our approach in the considered cases.

In the next step the normality assumption of the error distribution is replaced by a smooth distribution of Gaussian distributions in conjunction with a smooth random effects distribution. This work is in development. Up to now, the results show that a mixture of Gaussian distributions on the error distributions combines nicely with a Gaussian random effects distribution, but that, due to a large number of (simple) integrals to be evaluated a different numerical technique is needed.

#### **Joint Modelling of Multivariate Longitudinal Profiles**

In the context of longitudinal data or repeated measurements, research questions are often formulated which require joint modelling of multivariate response vectors measured repeatedly within the participating subjects. Examples include longitudinal studies where many indices are measured over time and where an overall assessment of timetrends is needed or where classification of subjects based on multivariate longitudinal profiles is of interest. Other examples can be found in clustered settings, e.g. a questionnaire measuring different concepts, each by a set of items.

There are a number of possible approaches to extend models for one repeatedly measured outcome to multivariate settings. We will focus on a random-effects approach, where a mixed model is assumed for each outcome separately, and where the joint model arises from assuming a joint (multivariate) distribution for all random effects. This model is very flexible in the sense that it does not assume the outcomes to be measured the same number of times, nor at the same time points. This approach also allows combining outcomes and/or models of a different nature: Continuous and categorical outcomes; linear, generalized linear, and non-linear models.

Nevertheless, there is also a disadvantage, especially when high-dimensional outcome vectors need to be analysed, which is the case in the examples we will discuss. A new pairwise model fitting approach has been developed which can circumvent this problem, and which is applicable whatever dimensionality of the problem (Fieuw and Verbeke (2004)). In this approach, the likelihoods of all pairwise models (each model involves two outcomes) are first maximised separately instead of the likelihood of the full multivariate model. In the second step, parameter estimates are obtained by averaging over all pairs. Borrowing ideas from the pseudolikelihood framework, standard errors can be calculated and inferences become readily available. The approach has been illustrated using examples with linear mixed and generalised linear mixed models, and the statistical properties of the estimators and associated standard errors are evaluated using simulation studies.

### **1.5.2 The investigation of mixture models as an alternative for approaching the random effects distribution**

The topics that are treated here are much related to those of **WP5**. The following topics were treated:

#### **Allowing for examiner's bias and variability in a logistic random effects model**

Mwalili, Lesaffre and Declerck (2005) developed a method to correct for the examiner's bias and variability when a gold standard and calibration data are available for an ordinal logistic regression model. The method was applied to the geographical distribution of the dmft-score (caries) of seven-year old children in Flanders. Further, Lesaffre, Mwalili and Declerck (2004) have explained the methodology in the previous paper to a dental audience. Further, the same methodology has been implemented to a Zero-inflated Poisson distribution for the dmft-score. Furthermore, in collaboration with Helmut Kuechenhoff, Mwalili and Lesaffre have developed the SIMEX (Simulation-Extrapolation) method for misclassification. A manuscript has been conditionally accepted by Biometrics (Kuechenhoff, Lesaffre and Mwalili (2004)). Two further developments are being explored. In both approaches we examine how the misclassification matrix used above to correct for misclassification can be estimated with higher precision. For this we look at the misclassification mechanism of the dmft-score as the sum of tooth-specific misclassifications. In a first approach, we still model a mouth-specific score but estimate the misclassification probabilities via the tooth-specific misclassification matrix. In a second approach we explore tooth-specific hierarchical models and estimate the misclassification matrix using the tooth type and the specific true status of the tooth. A manuscript is in preparation.

#### **Statistical framework for item response models**

De Boeck and Wilson (2004) published the book entitled 'Explanatory Item Response Models: A Generalized Linear and Nonlinear Approach' to the Springer series on Statistics for Social Science and Public Policy. It is an edited book, but written as a monograph, in collaboration with researchers involved in **WP4** (Geert Verbeke, Steffen Fieuw from KUL-2), **WP5** (Paul De Boeck, and others from KUL-1), and **WP6** (Geert Molenberghs from LUC) on one hand, and the educational measurement research group from UC Berkeley (Mark Wilson) on the other hand. As a consequence of this publication Paul De Boeck is invited to organize workshops at the meeting of the National Council for Measurement in Education (NCME) in Montreal (April 2005) and the

International Meeting of the Psychometric Society in Tilburg (July 2005). Authors from **WP4**, **WP5**, and **WP6** co-teach a graduate course on educational measurement at UC Berkeley.

### **Residual dependence**

A conditional approach for modeling residual dependency for mixed logistic models has been developed further, in the line of earlier work, but now for ordered-category data. In a first article, a model for partially-ordered-category data was formulated and compared with a similar one based on a partly marginal approach (Ip, Wang, De Boeck, and Meulders (2004)). In a second article, this model was adapted to handle residual dependency between rating scale data, and to investigate latent trajectories through a two-dimensional space of manifest categories (Meulders, Ip, De Boeck (2004)). The application describes trajectories through categories of intensity and frequency of emotions. As an alternative to the conditional approach, we have developed a dependency model for binary data based on Franks copula, so that the simple logistic form is retained in the marginal models. This line of modeling is a breakthrough in handling residual dependencies, because it allows for rather simple estimation and preserves a simple marginal form.

### **MIRID**

At KUL-1, De Boeck and Smits (2004) have further investigated models for covariate effects that are a function of other covariate effects (so-called MIRIDs). This has led to the formulation of a double-structure structural equation model. This is a model with multiple random effects for two modes of a three-mode data array. A crucial issue is the high dimensionality of the random effects. To deal with this problem, KUL-1 started to investigate Bayesian methods for high-dimensional models, and quasi-Monte-Carlo methods (Halton, Sobol) for high-dimensional numerical integration.

### **Counting process**

Tuerlinckx (2004) has proposed a multivariate counting process with positive dependencies for reaction times that can be approached as a random-effects model with independent non-homogeneous Poisson processes (conditional on the random effect).

#### **1.5.3 Extensions to interval-censored data**

The topics that are treated here are much related to those of **WP3**. We distinguish the following research topics:

#### **Modeling the emergence times using random effects models for interval, left and right censored data.**

The Signal Tandmobiel Study is a longitudinal dental study on about 4500 children. The KUL-2 team (Komárek and Lesaffre) fitted in collaboration with Tommi Härkänen several survival models with two random effects parameters (frailty parameter and birth of dentition parameter) for emergence the caries experience of the first permanent molars. This work is accepted in Biostatistics (and published early 2005), see Komárek, Lesaffre, Härkänen, Declerck and Virtanen (2005). The above analyses suggested examining a different approach, i.e. AFT-models (accelerated failure time models) with a complex error structure not necessarily assuming classical assumptions like normality. The KUL-2 team, in collaboration with Hilton (UCSF), has developed an AFT model

with a smooth error distribution being the mixture of Gaussian distributions. The manuscript has been accepted for *Journal of Computational and Graphical Statistics* (Komárek, Lesaffre, Hilton, 2005). An experience obtained during the work on interval-censored data was further summarized in the manuscript Lesaffre, Komárek and Declerck (2005) which was accepted by *Statistical Methods in Medical Research*. The manuscript also extends the model of Komárek, Lesaffre and Hilton (2005) by allowing also the scale parameter of the AFT model to depend on covariates. With respect to dental applications (modelling emergence and caries times), there is a need to include random effects in the model and to allow for doubly interval-censored data. Bayesian approach with its simulation based MCMC methodology avoiding explicit integration is an appealing alternative to the penalized maximum likelihood method used in the early AFT-developments of the KUL-2 team. A Bayesian model with normal random effects and an error structure being a classical (unpenalized) normal mixture with unknown number of mixture components was developed and submitted in October 2004 to *Statistica Sinica* (Komárek and Lesaffre (2004b)). An alternative to the classical mixture with unknown number of components is a penalized mixture with higher number of fixed components used already in the approach of Komárek, Lesaffre and Hilton (2005). Its Bayesian alternative used as an error and random effects distribution in the AFT model for multivariate doubly interval censored data is currently under development. Two manuscripts are anticipated in 2005, one of them in cooperation with the LUC group of Paul Janssen.

#### **Modeling jointly repeated measurements and survival times**

In the context of a clinical trial or an epidemiological study, there are often repeated measures on a risk factor available which have an impact on the survival response (e.g. serum cholesterol on survival). It is advantageous for the prediction of survival to model the repeated measurements model and the survival model jointly. This is often done by assuming conditional independence of the repeated measurement and the survival outcome conditional on some well-chosen random effects. The IAP doctoral student Dora Kocmanova explored different Bayesian techniques. Dora Kocmanova has been awarded a Marie-Curie Training Site Scheme Grant in the University of Lancaster (UK). Under the supervision of Peter Diggle and Rob Henderson she explored the repeated measurements-survival techniques from Jan 1, 2004 until 30 June, 2004. Her thesis topic is the development of models for jointly modeling an interval censored response with repeated measures taken on continuous/discrete covariates. Further, the existence of latent classes in the data will also be examined, the latter is important for modeling HIV data.

#### **Modeling multivariate interval-censored emergence times**

A GEE-method for modeling multivariate interval-censored data has been developed, and has been applied to the emergence times of the permanent teeth as recorded in the Signal Tandmoebel Study. Furthermore an improvement to the current methodology to calculate the bivariate non-parametric estimate of a survival function for interval-censored data has been proposed by Bogaerts and Lesaffre (2004a). Lesaffre and Bogaerts (2004) developed also a smooth bivariate estimate of the survival function, using a penalized likelihood approach. This methodology can be employed to have an improved estimate of association measures in bivariate survival models. In this respect, Lesaffre and Bogaerts developed an estimator for Kendall's tau using the bivariate smooth estimate of the survival distribution. The performance of this estimator in small sample



sizes has been explored and good results were obtained. Also another association measure, the Spearman correlation coefficient, can be easily estimated using the same technique. A manuscript describing this methodology had been submitted to *Biometrics* but was rejected and will be soon resubmitted to another journal after adaptation to some comments of the referees. In contrast to the first penalty, which works conditionally on both dimensions and has two smoothing parameters, a more natural bivariate penalty with only 1 smoothing parameter has also been investigated. The results of both penalties are to be compared in the near future as the simulations for the second penalty are almost finalized. Finally, we are investigating how covariates can be included in the model.

Cecere and Lesaffre (KUL-2) are developing a Bayesian model to estimate the correlation in a multivariate survival model when the survival times are interval-censored. At this stage only normally distributed interval-censored data are explored. The method has been applied to the emergence times of some teeth from the children sampled in the Signal Tandmobiel Study. Covariates are possible in this model and this allows us to see patterns in the emergence times (and associations between them) possibly depending on gender, dietary behavior and the geographical location in Flanders. Currently, we are studying the extension to general copula models and for this we look for collaboration with groups from other working packages (see **WP1** (Section 1.2.9), and **WP6**).

#### 1.5.4 Modeling compliance data

Rizopoulos, Tsonaka and Lesaffre (KUL-2) developed frequentist methods to analyze U- and J-shaped cross-sectional data. This kind of data occurs frequently in compliance studies as the percentage of days that a patient takes his/her drug correctly. The idea behind the models is to replace the observed proportions by a latent score which has on a transformed scale (logit-scale) approximately a normal distribution. This allows to employ classical (random effects) models (in a Bayesian context) to analyze these data. A manuscript (Lesaffre, Rizopoulos and Tsonaka (2004)) has been submitted to *Biostatistics* and is under revision now. In a second step we have explored the determination of the power and sample size for this model. A manuscript has been submitted to *Biometrics*, see Tsonaka, Rizopoulos and Lesaffre (2005). A third manuscript on the use of this model in a repeated measurements context is in preparation.

#### 1.5.5 Mixed models in actuarial sciences

Mixed models are used in actuarial science to account for some (residual) heterogeneity among the risks passed to the insurance company. A posteriori ratemaking techniques are used to correct the premium amount for this heterogeneity. In that respect, bonus-malus systems and credibility theory are very efficient techniques to refine the a priori risk classification. The papers by Denuit and Lang (2004), Pitrebois, Denuit and Walhin (2004) and Purcaru, Guillén and Denuit (2004) propose new methods to reevaluate the future premiums taking into account the past claims history.

### 1.5.6 Functional mixed models

This section combines features from **WP1** and **WP4**. Functional mixed-effects models are very useful in analyzing functional data. Antoniadis and Sapatinas (2004) consider a general functional mixed-effects model that inherits the flexibility of linear mixed-effects models in handling complex designs and correlation structures. Wavelet decomposition approaches are used to model both fixed-effects and random-effects in the same functional space. This helps in interpreting the resulting model as a functional data model since it does not contradict the intuition that, if each outcome is a curve, which is the basic unit in functional data analysis, then the population-average curve and the subject-specific curves should have the same smoothness property (i.e., they should lie in the same functional space).

## 1.6 Work package 5: Classification and mixture models

The progress on the various primary objectives as described in the research proposal is briefly described below.

### 1.6.1 Studying specific types of mixture models

The specific types of mixture models we have concentrated on are of four types.

#### PMD models

An important achievement is the hierarchical extension of the model, with hyperparameters referring to the distribution of the component probabilities (Meulders, De Boeck, Van Mechelen, and Gelman (2004)). The hyperparameters can inform us about the categorization tendency in processes underlying the perception of emotions.

#### Item response models

Item response models (IRT) are models for repeated binary and ordered-category data with a logit or probit link. Mixture distributions can be used for the random effects.

1. An important topic is the differentiation between dimensional and categorical structures. On the one hand, we have investigated the relationship between a model with heterogeneous mixture components and a continuous multidimensional model (Rijmen and De Boeck (2004)). On the other hand, we have developed a method deemed Dimcat to differentiate between latent categorical and dimensional structures (De Boeck, Wilson and Acton (2005)). This work is published in the top journal of psychology (impact factor > 8). It is a serious competitor for the dominant taxometric approach to detect discontinuities in psychological phenomena (development, psychiatric syndroms, attitudes and opinions). The ideas for this model stem from our work on differential item functioning based on a generalized linear mixed model approach (Van den Noortgate and De Boeck (2004)).

2. We have developed and investigated models for change. First, one of the achievements is a model for transitions between heterogeneous mixture components, each representing a developmental stage with individual differences (Rijmen, De Boeck and van der Maas (2004)). Because of the equivalence of an IRT model (with a discrimination parameter) with a diffusion model of a certain type (Tuerlinckx and De Boeck (2004)), also transitions between qualitatively different

diffusion models can be dealt with in this way. Second, models for change points within a series of repeated measures are studied. Either the change point is fixed, as applied by Smits and De Boeck (2004), or it is random, and it can be approached with a mixture model, which is what we started to explore. Third, also the trajectories of change can be a subject of mixture modeling, as in applications by De Fraine, Van Damme, and Onghena (2004), and by Prinzie, Onghena, and Hellinckx (2004). Growth mixture models were applied on children's aggressive and delinquent problem behavior, and on data regarding adolescents' academic self-concept, in order to distinguish between several developmental pathways.

### **Models with an errant process**

It often occurs in psychological studies that not all data are generated through the assumed process, but instead through a minority process that is considered as an errant process. Guessing is an example. This process can be captured in a mixture component that is differentiated from the dominant normal process. We have expanded a Rasch model with random effects for the success probabilities of the errant process component (San Martin, Del Pino, and De Boeck (2004)). This work has led us to model how at the lower extreme of the ability scale, less able persons can obtain better results. This phenomenon can be captured if the success probabilities of the errant process component are a negative function of the common underlying latent variable (the ability).

### **Models with smoothing approaches**

Also smoothing approaches are developed and investigated for their use in the context of mixture modeling. First, at the KUL-2 team, Ghidry, Lesaffre and Eilers (2004) have developed, in collaboration with Paul Eilers (University of Leiden) a mixed model with a smooth random effects distribution. For instance, the model allows a long-tailed distribution for the random effects, as well as a mixture of normal distributions. The smoothing approach will reveal the underlying mixture structure, though without explicitly detecting the components of the mixture. Penalized likelihood maximisation delivers the estimated fixed effects and the smoothing coefficients of the random effects distribution. A similar approach was developed for survival models with left-, right and interval censoring: Komárek, Lesaffre and Hilton (2005) developed an AFT model with a smooth error distribution (see also **WP3**, Section 1.4.2). Results show that this model nicely reveals the mixture structure of the error distribution, if present. Bogaerts and Lesaffre (2004b) have further developed this approach for fitting a bivariate survival model. A Bayesian model with normal random effects and an error structure being a classical (unpenalized) normal mixture with unknown number of mixture components was developed and submitted in October 2004 to Statistica Sinica (Komárek and Lesaffre (2004b)). An alternative to the classical mixture with unknown number of components is a penalized mixture with higher number of fixed components used already in the approach of Komárek, Lesaffre and Hilton (2005). Its Bayesian alternative used as an error and random effects distribution in the AFT model for multivariate doubly interval censored data is currently under development (see also **WP3**).

Second, in collaboration with the UCL team, Taoufik Bouezmarni (now a IAP-postdoc of the KUL-1 team, but previously at the UCL) investigates the possibilities of kernel smoothing to approach the underlying distribution of a random effect. The idea is to begin with the estimation of homogeneous mixture components following a jackknife procedure, such that multiple locations

along the scale are obtained, one for each run. Starting from the resulting parameter estimates, the kernel method can be used in a second step for a nonparametric approach to the random effects distribution. This work forms a natural bridge between **WP1** and WP5.

### **1.6.2 Investigating methods to decide on the number and type of components**

A comprehensive approach for model checking in models with missing or latent data has been developed by members of the KUL-1 and KUL-2 team (Gelman, Van Mechelen, Verbeke, Heitjan, and Meulders (2004)). This approach, proposed within a Bayesian framework, extends the multiple imputation principle from the domain of model estimation to that of model checking, and especially relies on graphical tools. It includes the case of mixture models (which involve latent data) and checking problems with regard to the number and type of mixture components as a special case. The number of clusters issue has further been given special attention within the study of classification techniques other than mixtures (see below).

### **1.6.3 Classification techniques other than mixtures**

We studied a one-mode additive clustering model for two-way two-mode data. The mathematical characteristics of this model (including the special case of a model with nested clusters) have been investigated and a novel algorithm to estimate it has been developed (Depril and Van Mechelen (2004)). Also, we finalized work on a conjunctive model for picky any/n data that represents choice behavior in terms of clusters in a low dimensional space (defined by a set of ordinal variables with a prespecified number of categories) (Leenen and Van Mechelen (2004)).

A considerable amount of research has further been devoted to the simultaneous clustering of several modes in case of two-mode or multimode data, also in collaboration with the RWTH-Aachen team. First of all, a comprehensive taxonomy of simultaneous two-mode clustering methods has been developed by the Leuven and Aachen teams, implying a set of novel structuring principles for this fairly heterogeneous domain (Van Mechelen, Bock, and De Boeck (2004a,b)). The RWTH-Aachen team was also a co-editor of a volume on multivariate analysis, including clustering (Chiodi, Mineo, and Bock (2004)), and several types of clustering have been investigated (Bock (2003a); De Carvalho, Brito, Bock (2004)). This latter work refers mainly to the analysis of ‘symbolic data’ of the multivariate interval type, a quite general scenario of ‘complex data’.

The design and mathematical investigation of a ‘convexity-based’ clustering criterion for analyzing heterogeneity in data sets is an important topic of research in the Aachen group. The results of a theoretical study, a description of algorithms, and applications can be found in Bock (2003b) from RWTH. This kind of clustering is directly suited to the simultaneous clustering of the rows and columns of a contingency table such that the resulting row and column partitions are ‘maximally related’ to each other. Furthermore, the work of Hans Bock on proximity measures to start from in clustering is recognized with an article in the Encyclopedia Statistics in Behavioral Science (Bock (2004)). In addition, the Aachen group has also considered optimization problems in the context of portfolio selection in cases where the literature has so far even not yet proved the existence of an optimum strategy. In his dissertation, Beutner (2004) proposed, and investigated, conditions that guarantee the existence of an optimum strategy by reducing this problem to the

question if, for two closed subspaces (convex sets)  $A, B$  of a Hilbert space, the sum set  $A+B$  is closed as well.

The KUL-1 group has further investigated simultaneous partitioning approaches for multi-mode data. In particular, for the three-mode case a novel algorithm has been developed, and its performance has been evaluated and compared with that of two alternative algorithms in a comprehensive simulation study (Schepers and Van Mechelen (2004)). The KUL-1 group finally has extensively investigated novel extensions of the hierarchical classes (HICLAS) family, which is a comprehensive family of simultaneous classification models for multimode data. The research in this regard has proceeded along five different lines: (1) the development of models and associated algorithms for two-way two-mode real-valued data; up to now, the hierarchical classes family was limited to binary and categorical (rating-valued) data; the novel extension of the family implies a considerable enlargement of its scope. (2) The development of new algorithmic variants (based on new types of starting configurations and simulated annealing-type algorithms) to deal with enduring local minima problems, along with the evaluation of those algorithms in extensive simulation studies (Ceulemans and Van Mechelen (2004a)). (3) The addition of confidence information to point estimates of parameter vectors through nonparametric bootstrap methods and through MCMC estimation procedures of minimal stochastic extensions of deterministic HICLAS models (Leenen, Van Mechelen, Gelman, and De Knop (2004)); one may note that the latter approach also paves the way for novel types of model checking. (4) The development of new tools for model selection, especially with regard to the modeling of three-mode data; more in particular, a novel convex hull based automated model selection procedure has been developed for three-mode HICLAS models, which appeared to perform excellently in simulation studies (Ceulemans and Van Mechelen (2004b)). (5) The development of novel three-mode HICLAS models that do not reduce all three modes of the data (Ceulemans and Van Mechelen (2004c)), that include constraints based on theoretical or empirical prior knowledge (Ceulemans, Van Mechelen and Kuppens (2004)), and that represent heterogeneity in sequential, chain-type processes.

#### 1.6.4 Specific cross-links with other work packages

An important part of the work done by the KUL-1 group concerns generalized linear and nonlinear mixed models with a logit link, a topic that is directly related to **WP4** (mixed models) and is also of relevance to **WP6** (latent variables).

#### 1.6.5 Methodological problems

The problem of heterogeneity is an important one for a large variety of models. We have investigated six methods to detect and locate heterogeneity in logistic regression models (individual differences in the intercept and slopes). As mentioned earlier a GEE2 approach turned out to be very successful (Balazs, Hidegkuti, and De Boeck (2004)), but in an additional study which is now included in the revised and resubmitted manuscript, they found that a nonparametric method (DETECT) can be adapted to yield good results as well. Further, Katalin Balazs is investigating a method to combine an additive clustering method (ADCLUS) with a GEE2 approach to deal with multidimensional heterogeneity.

### 1.6.6 Functional classification techniques

This section combines features from **WP1** and **WP5**. In Le Borgne, Guérin and Antoniadis (2004) (UJF), Independent Component Analysis (ICA) is used to compute features extracted from natural images. The use of ICA is justified in the context of classification of natural images for two reasons. On the one hand the model of image suggests that the underlying statistical principles may be the same as those that determine the structure of the visual cortex. As a consequence, the filters that ICA produces are adapted to the statistics of natural images. On the other hand, they adopt a nonparametric approach that requires density estimation in many dimensions, and independence between features appears as a solution to overthrow the ‘curse of dimensionality’.

Bugli, Lambert, Boulanger, Ledent, Pereira and Nardone (2004) describe statistical tools to detect and quantify the effect of drugs on the brain from electroencephalograms. ICA is first used to detect and to remove artefacts from the observed signals, and then again to reduce the dimension of the problem. They show that eight components can reconstruct more than eighty percent of the data recorded using twenty-eight electrodes. Some of these components correspond to an important characteristic of the signal named event-related potential. Two of them are particularly interesting because they decompose the P300 peak.

Antoniadis, Bigot and von Sachs (2005) (from UJF and UCL) are working on a research project on the development of multiscale and multigranular functional classification methods to recognize particular features in multispectral MRI brain images.

The work on classification trees has been pursued and I. Demacq (advisor L. Simar) from UCL has defended her thesis in 2004. In Demacq and Simar (2005), an exact algorithm based on hyper-rectangular partitioning trees has been implemented. It has optimal properties but the computing complexity limits its applicability. An approximated (faster) algorithm, based on the quantiles for building the splitting rules has been implemented. Its performances are compared with more classical classification algorithms through some classical data sets and it appears that, for these data sets, this new classifier performs much better.

## 1.7 Work package 6: Incompleteness and latent variables

The work on incomplete data can be divided, broadly, in complex modeling approaches for incomplete data and sensitivity analysis tools.

Lipsitz, Molenberghs, Fitzmaurice and Ibrahim (2004) present methodology to estimate parameters for models where missingness is missing not at random. They use a so-called protective estimator, where the missingness mechanism need not be addressed explicitly, yet valid estimators can be obtained.

Methods to properly deal with incomplete longitudinal clinical trial data have been proposed by Molenberghs, Thijs, Jansen, Beunckens, Kenward, Mallinckrodt, and Carroll (2004). Mallinckrodt, Watkin, Molenberghs, and Carroll (2004) advocated the use of proper modeling approaches, rather than ad hoc methods, for the primary analysis of longitudinal clinical trials. The type I error rate when likelihood-based repeated measures analyses are used in incomplete longitudinal data was studied in Mallinckrodt, Kaiser, Watkin, Molenberghs, and Carroll (2004).

The book by Dmitrienko, Offen, Faries, Christy Chuang-Stein, and Molenberghs (2004) contains

a large chapter on the proper treatment of incomplete clinical trial data, with an emphasis on models that can be implemented using the SAS Software system.

Molenberghs and Verbeke (2004) wrote a chapter in the book by De Boeck and Wilson (2004) on repeated measures, mixed models (interaction with WP4), and incomplete data. The book as a whole and this chapter in particular aims to bridge the gap between the psychometric and biostatistical research communities.

Verbeke (KUL-2) and Molenberghs (LUC) actively disseminate incomplete data and longitudinal methodology through a variety of short courses, taught in various continents. They will publish a Springer book (2005) on longitudinal and incomplete data methods for non-Gaussian data.

Interaction with **WP3** : Tibaldi, Molenberghs, Burzykowski, and Geys (2004) propose a pseudo-likelihood based estimation for a marginal multivariate survival model, based on the use of copulas. Tibaldi, Torres Barbosa, and Molenberghs (2004) proposed a marginal method to assess associations between time-to-event responses in a pilot cancer clinical trial.

The remaining of this section has strong links with **WP4**. In the context of surrogate marker evaluation, Molenberghs, Burzykowski, Alonso, Geys, and Cortiñas have proposed a mixed-model based methodology, which incorporates both multivariate outcomes, surrogate and true outcomes in clinical trials, and hierarchical data, resulting from several trials. Molenberghs, Burzykowski, Alonso, and Buyse (2004) provided a general overview. A reflection on the pros and cons of the meta-analytic paradigm, in contrast to the classical approaches, was presented in Alonso, Molenberghs, Burzykowski, Renard, Geys, Shkedy, Tibaldi, Cortiñas, and Buyse (2004). In such hierarchical data, one can never be sure all relevant hierarchical levels have been taken into account. This problem occurs in the context of surrogate markers, but more broadly as well. Cortiñas, Molenberghs, Burzykowski, Shkedy, and Renard (2004) made a detailed study of this topic. Various therapeutic areas require specific methodology. A case study in oncology, based on an ordinal or binary surrogate for a time-to-event endpoint was treated by Burzykowski, Molenberghs, and Buyse (2004). When the meta-analytic framework is further expanded by allowing the surrogate and true endpoints to be measured repeatedly, specific methodology is needed. Such methods were proposed in Alonso, Geys, Molenberghs, Kenward and Vangeneugden (2004).

Molenberghs and Verbeke (2004) indicate how complex models, including but not restricted to mixed models, should be examined for their meaningfulness. Especially generalized linear and non-linear mixed models are prone to a number of subtle complexities that may easily be overlooked by the practitioner, and a fundamental reflection is therefore at hand.

In the context of hierarchical binary data, pseudo-likelihood ideas based on a probit model formulation are used by Renard, Molenberghs, and Geys (2004).

Reliability estimation is typically done in the context of psychometric or sociological studies. However, it is also of relevance in clinical and epidemiological studies, e.g., in the context of psychiatric clinical trials, quality of life assessment, etc. To allow using clinical trial data for this purpose, Vangeneugden, Laenen, Geys, Renard, and Molenberghs (2004) proposed a linear mixed model approach to assess reliability. These methods were further expanded to so-called generalizability theory in Vangeneugden, Laenen, Geys, Renard, and Molenberghs (2004).

## 2 Network activities

### 2.1 Web site

A new web site has been set up in 2004, which is meant to be more user friendly both for the webmaster as for any person consulting the site. The creation of the new web page has been in charge of Eric Lecoutre, a computer expert of the Institute of statistics at UCL.

All activities of the IAP-statistics network can be followed very closely from our web site. The address of the web site is <http://www.stat.ucl.ac.be/IAP>. The web site contains e.g. the following information :

- Our logo
- Call for applications
- Description of the project
- List of scientific personnel working under IAP project
- Downloadable member list
- Research activities (workshops, seminars, short courses,...)
- Downloadable technical reports, list of publications and list of books written by members of the network
- Annual reports and reports of scientific meetings
- Contact details

### 2.2 Technical reports and published papers

Two publication series, available via the web site, report on scientific results obtained within the IAP-statistics network: the technical report series and the list of publications. The IAP-statistics technical report series groups all papers written in the IAP-statistics network. Each paper in this series has been submitted for publication to an international journal. Once a paper has been accepted for publication in an international journal and has been printed, we list it into the IAP-statistics list of publications. For the IAP-statistics technical report series we put for each paper the pdf file of the complete paper on our web site. For the IAP-statistics list of publications, we provide on the web site a list of references of the published papers.

### 2.3 Scientific meetings

#### 2.3.1 Workshops

On 30 November 2004 a half day meeting on 'Frailty models' has been organized in Leuven by the KUL-2 and LUC partner. Four researchers (two from inside and two from outside the network) have given invited presentations at this meeting. More information can be found on the web page <http://www.stat.ucl.ac.be/IAP/download/meeting30nov04.pdf>.



An international workshop was planned to be organized by the ULB partner during 2004. A keynote speaker of this meeting would have been Steve Portnoy (University of Illinois, USA) who planned to visit the ULB in 2004. Due to changes in the dates of his visit, the meeting had to be rescheduled and will take place on 12-13 May 2005. The workshop will be organized on the general theme of ‘Asymptotics’ and will host speakers from inside and outside the network.

In 2005, a second international meeting will be organized by the KUL-1 partner, probably at the end of September. The theme of this workshop will be ‘How to deal with heterogeneity’. In this regard, we want to use a broad concept of heterogeneity that covers aspects from many of the work packages like, for instance, nonstationarity in time series analysis, random effects as included in mixed and survival (frailty) models, as well as the presence of different subpopulations as assumed in mixture models. The presentations at this workshop will consist of joint contributions on the theme of heterogeneity, each from at least two network partners. The joint papers can be the result of collaborative work in the form of joint theoretical work, of data from one partner that are analyzed by another partner, or of any other form.

In 2006 the coordinator at UCL will be the host for a fifth and final workshop, during which the achievements of the past years will be summarized. Members of the KUL-2 and UCL team plan to organize a one-day meeting around the topic of interval censored data near the end of 2005. More details will follow.

### **2.3.2 Special activities by young researchers**

After the success of the ‘First Young Researchers Day’, organized by the PhD students of the Institute of statistics of UCL in May 2003, a ‘Second Young Researchers Day’ was organized on April 30 2004, again at UCL. Among the main objectives of this day was the exchange information around common research interests of PhD students in statistics in Belgium. The one-day meeting was organized around the topic of ‘Many explanatory variables ? A challenge for regression modelling’ and was attended by 86 participants. The one-day meeting was also supported by the Graduate School of the Institute of Statistics (UCL), the IMS (Institute of Mathematical Statistics) and the FNRS. Such a meeting is a real encouragement for PhD students, and enforces communication between them. For the programme of this meeting, and other details, see the web site <http://www.stat.ucl.ac.be/YRD/YRD2/mainpage.html>. A third such meeting is foreseen for Fall 2005.

## **2.4 Organization of the network : administrative meeting**

The annual administrative meeting with a representative from the federal office OSTC-Brussels (Ms Lejour) took place on 17 December 2004 at UCL. Two members from the follow-up committee (Prof. T. Snijders and Prof. M. Delecroix), the promoters of the network, the coordinators L. Simar and I. Van Keilegom and Ms C. Denayer (head of the administration at the Institute of statistics at UCL) and P. De Boeck, Y. Van Mechelen, N. Veraverbeke, M. Hallin and H. Bock participated to the meeting. Other promoters of the network joined us in the afternoon for the scientific part of the meeting (two seminars, see below).

## 2.5 Collaborations, working groups and seminars

### 2.5.1 Collaborations

The number of scientific collaborations within the network continues to grow, as can be seen from the list of technical reports and publications. During 2004 a number of members has visited one of the European partners in order to collaborate efficiently on ongoing projects (e.g. Irène Gijbels and Rainer von Sachs visited Anestis Antoniadis (UJF) for 1 week and 2 weeks respectively). That there exists a vivid collaboration spirit among the members of the network can also be seen from the list of short courses outlined in Subsection 2.6 (some of these short courses are taught or organized jointly by several teams of the network) and from the list of postdoc positions financed by the IAP outlined in Subsection 2.7 (some of the PhD students of the network get postdoc/professor positions at other universities of the network). This will help establishing an active research network on the long run.

Below, we mention a few examples of ongoing collaborations between members of different teams of the network.

#### *Frailty models and inference*

To compare the different models available to study heterogeneity in multivariate survival data a working group has been installed. The network members are Lesaffre and Komárek (KUL-2), Lambert (UCL) and Janssen and Massonnet (LUC). Other participants include Legrand and Sylvester (EORTC) and Duchateau (UGent). The coordinating team is the LUC-team (Janssen).

#### *Goodness-of-fit problems*

A growing number of researchers from several teams are working on goodness-of-fit problems for the shape of a regression function (Claeskens (UCL, Leuven), Aerts (LUC) and Van Keilegom (UCL)). Gerda Claeskens is currently writing a book on this topic (joint with Nils Hjort (University of Oslo)).

#### *Interval censored data*

Researchers from KUL-2 (Lesaffre and some of his students) and from UCL (Van Keilegom) started collaboration on interval censored data problems. Expertise on theoretical issues (UCL) and on applied issues (KUL-2) are here nicely combined and lead to better insight from both sights. A workshop on this topic will be organized jointly by both research groups at the end of 2005.

### 2.5.2 Working groups

The frailty working group continued its series of meetings. In 2004 the group met at several occasions to discuss research progress and for a series of seminars mainly to introduce relevant methodology to new PhD students. Participants were Burzykowski (LUC), Cortiñas (LUC), Duchateau (UGent), Janssen (LUC), Legrand (EORTC), Massonnet (LUC), Sylvester (EORTC) and for some of the meetings Ampe (UGent), Goethals (UGent), Ducroq (INRA), Lesaffre (KUL), Komarek (KUL) and Lambert (UCL).

At LUC, an active research group on surrogate marker evaluation (Burzykowski and Molenberghs) is based, with participation of colleagues from various other groups, including pharmaceutical companies (Vangeneugden, Virco-Tibotec; Tibaldi and Renard, Eli Lilly) and contract organizations (Buyse International Drug Development Institute). The team will publish a Springer book on the topic in 2005. Likewise, there is an active research group on incomplete data and sensitivity analysis (Molenberghs and Aerts) with collaboration from Leuven (Verbeke), and participation from pharmaceutical companies (Mallinckrodt, Eli Lilly US), Texas A&M University (Carroll), and the London School of Hygiene and Tropical Medicine (Kenward).

### 2.5.3 Seminars

Each of the participating partners organizes on a regular basis statistics seminars at their universities (on the average 3 research seminars per week during the academic year). Announcements of these seminars are sent out to most of the Belgian statisticians, including those participating in the network.

Apart from the regular statistics seminars at the universities involved, several other seminars have been organized in the IAP-statistics network, like e.g.

- Seok-Oh Jeong (UCL), ‘Linearly interpolated FDH estimator for nonconvex frontiers’, 29 October 2004, UCL
- Sébastien Van Bellegem (UCL), ‘Semiparametric estimation by model selection for locally stationary processes’, 17 December 2004, UCL
- Tom Snijders (University of Groningen, The Netherlands), ‘Frequentist MCMC estimation methods for multilevel logistic regression’, 17 December 2004, UCL

In addition, a number of seminars have been given by members of the network at other universities of the network, in order to stay informed of the research activities of the other partners.

## 2.6 Short courses

Several short (intensive) courses have been organized within the framework of the IAP-statistics network. These courses were intended for all members of the network, and in particular (but not exclusively) for the PhD-students. The announcements were each time sent out to all members and posted on the web site. No (or reduced) registration fees were required for IAP-members.

A list of the short courses organized during the working year 2004 is given below.

- An intensive course of 7.5 hours on ‘The power of penalties’ given by Professor Paul Eilers (Leiden university) in December 2003 and February 2004 at UCL.
- An intensive course of 7.5 hours on ‘Frontier estimation and the use of bootstrap methods’ given by Professor Paul Wilson (Texas University, Austin, USA) in February-March 2004 at UCL.
- An intensive course of 7.5 hours on ‘Design of experiments’ by Professor Holger Dette (University of Bochum, Germany) in February-March 2004 at UCL.

- An intensive course of 15 hours on ‘Smoothing techniques and nonparametric testing’ given by Professor Gerda Claeskens (UCL, ‘terugkeermantat’/‘mandat de retour’ from the DWTC/SSTC) in Spring 2004 at UCL.
- An intensive course of 7.5 hours on ‘Advanced bootstrap methods’ by Professor Peter Hall (Australian National University) in November-December 2004 at UCL.

For 2005, the following (incomplete) list of courses is planned :

- Series of lectures (11 lectures of 2 hours) given by Steve Portnoy (University of Illinois, USA), receiver of the ‘Chaire Francqui interuniversitaire au titre étranger’ (shared between ULB and KUL) in Spring 2005 at ULB and KUL.
- An intensive course of 7.5 hours on ‘Statistical denoising of images, with applications in astronomy’ by Véronique Delouille (Royal Observatory of Belgium) in Spring 2005 at UCL.
- An intensive course of 2 days on ‘Mixed Models and Incomplete Data’ by Geert Verbeke (KUL-2) and Geert Molenberghs (LUC) on 31 March and 1 April 2005 at UCB Pharma (Braine l’Alleud).
- An intensive course of 2 days on ‘Introduction to cluster analysis’ by Hans-Hermann Bock (RWTH) and Iven Van Mechelen (KUL-1) on 2-3 May 2005 at KUL-1.
- An intensive course of 2 days on ‘Introduction to bayesian inference in biomedical applications’ by Nicky Best (Imperial College, London) on 3-4 May 2005 at UCL.
- An intensive course of 7.5 hours on ‘Nonparametric regression techniques’ by Byeong Park (Seoul National University, South Korea) in Spring 2005 at UCL.
- An intensive course of 7.5 hours on ‘Empirical likelihood and its applications in nonparametric curve estimation and testing’ by Song Chen (Iowa State University, USA) in June-July 2005 at UCL.

The short courses organized by the UCL were also part of the doctoral programme of the graduate school in statistics at UCL.

## 2.7 Postdoctoral researchers and return grants

Li Chun Wang (PhD from Chinese Academy of Sciences, Beijing) was appointed as postdoctoral researcher of the network at LUC from 15 February 2004 to 15 February 2005. His main area of research is Bayes and empirical Bayes methods. During his stay he cooperated with Noël Veraverbeke and he gave research seminars at LUC and at UCL. Valentin Zelenyuk (Ukrainian, PhD from State University of Oregon) is now a post-doctoral researcher at the UCL, while Seok-Oh Jeong (PhD from University of South Korea) and Jérémie Bigot (PhD from UJF) ended their postdoc position at UCL funded by the IAP network in the course of 2004. Finally, Taoufik Bouezmarni (Maroccan, PhD from UCL) was appointed as postdoctoral researcher at KUL-1.

Gerda Claeskens, who obtained a ‘Return grant’ from the OSTC-Brussels in 2003 left the UCL in September 2004 and has now obtained an assistant professor position in Leuven.

### 3 Technical reports and publications

Below we provide in each of the subsections two lists of scientific works related to the IAP-statistics network:

#### A. List of Technical Reports:

This list contains all Technical Reports that have been written in 2004, and have been submitted for publication to an international journal. These reports are also available on our web site and the number listed refers to the electronic IAP-Statistics Technical report Series.

#### B. List of Publications:

This list contains all publications in international journals (with refereeing system), including also papers that are accepted for publication and are ‘in press’. This list also includes papers that have been published in Proceedings and have undergone a peer review (i.e. full length papers). See also the IAP-Statistics Reprints Series on our web site.

The list of Technical Reports is included since it allows us to provide a more complete overview of the achieved research results.

### 3.1 List of publications per research unit/partner

#### 3.1.1 Université catholique de Louvain, UCL partner

##### A. List of Technical Reports

Antoniadis, A. and Bigot, J. (2004). Poisson inverse problems, IAP-statistics Technical Report Series TR # 0425.

Antoniadis, A., Bigot, J. and Gijbels, I. (2005). Penalized wavelet monotone regression (in preparation).

Antoniadis, A., Bigot J. and von Sachs, R. (2005). Statistical analysis and characterization of brain response images (in preparation).

Antoniadis, A., Gijbels, I. and Nikolova, M.(2005). Penalized likelihood regression for generalized linear models with nonquadratic penalties (in preparation).

Badin, L. and Simar, L. (2004). A bias corrected nonparametric envelopment estimator of frontiers. IAP-statistics Technical Report Series TR # 0410.

Bonacorsi, A., Daraio, C. and Simar L. (2004). Advanced indicators of productivity of universities: an application of robust nonparametric methods to Italian data. IAP-statistics Technical Report Series TR # 0427.

Bugli, C. and Lambert, P. (2004). Functional ANOVA with random functional effects: an application to event related-potentials modelling for electroencephalograms analysis. IAP-statistics Technical Report Series TR # 0431.

- Bugli, C., Lambert, P., Boulanger B., Ledent E., Pereira, A. and Nardone, P. (2004). Statistical analysis of electroencephalograms: independent component analysis of event-related potentials. IAP-statistics Technical Report Series TR # 0468.
- Cao, R. and Van Keilegom, I. (2004). Empirical likelihood tests for two-sample problems via nonparametric density estimation. IAP-statistics Technical Report Series TR # 0429.
- Daouia, A. and Simar, L. (2004). Nonparametric efficiency analysis: a multivariate conditional quantile approach. IAP-statistics Technical Report Series TR # 0419.
- Daraio, C. and Simar, L. (2004). A Robust Nonparametric Approach to Evaluate and Explain the Performance of Mutual Funds. IAP-statistics Technical Report Series TR # 0413.
- Daraio, C. and Simar, L. (2005b). Conditional nonparametric frontier models for convex and non convex technologies: a unifying approach. IAP-statistics Technical Report Series TR # 0503.
- Delouille, V., Jansen, M. and von Sachs, R. (2004). Second generation wavelet methods for denoising of irregularly spaced data in two dimensions. Revised version of IAP-statistics Technical Report Series TR # 0303.
- Denuit, M., Purcaru, O. and Van Keilegom, I. (2004). Bivariate Archimedean copula modelling for loss-ALAE data in non-life insurance. IAP-statistics Technical Report Series TR # 0422.
- Färe, R. and Zelenyuk, V. (2004). On aggregation of scale elasticities across firms. IAP-statistics Technical Report Series TR # 0433.
- Geenens, G. and Delecroix, M. (2005). A survey about single-index models theory. IAP-statistics Technical Report Series TR # 0511.
- Gijbels, I., Lambert, A. and Qiu, P. (2004). Jump-preserving regression and smoothing using local linear fitting: a compromise. IAP-statistics Technical Report Series TR # 0401.
- Guyot C., François, N., Hug, B., Callemien, D., Govaerts, B. and Collin, S. (2004). Beer asstringency assessed by quantitative descriptive analysis and time intensity : influence of pH, oxygen and accelerated staling. IAP-statistics Technical Report Series TR # 0469.
- Henderson, D. J. and Zelenyuk, V. (2004). Testing for catching-up: statistical analysis of DEA efficiency estimates. IAP-statistics Technical Report Series TR # 0434.
- Hjort, N.L., McKeague, I.W. and Van Keilegom, I. (2004). Extending the scope of empirical likelihood. IAP-statistics Technical Report Series TR # 0415.
- Jeong, S.-O. (2004). Asymptotic distribution of DEA efficiency scores. IAP-statistics Technical Report Series TR # 0424.
- Jeong, S.-O. and Park, B.U. (2004). Limit distribution of convex-hull estimators of boundaries. IAP-statistics Technical Report Series TR # 0423.

- Jeong, S.O. and Simar, L. (2005). Linearly interpolated FDH efficiency score for nonconvex frontiers. IAP-statistics Technical Report Series TR # 0502.
- Kumbhakar, S.C., Park, B.U., Simar, L. and Tsionas, E.G. (2004). Nonparametric stochastic frontiers: a local likelihood approach. IAP-statistics Technical Report Series TR # 0418.
- Lambert, P. and Eilers, P.H.C. (2004). ‘Bayesian survival models with smooth time-varying coefficients using penalized Poisson regression’. IAP-statistics Technical Report Series TR # 0435.
- Mouchart, M. and Vandresse, M. (2004). Bargaining powers and market segmentation in freight transport. IAP-statistics Technical Report Series TR # 0421.
- Nguti, R., Claeskens, G. and Janssen, P. (2004). One-sided tests in shared frailty models. IAP-statistics Technical Report Series TR # 0436.
- Pardo Fernández, J.C., Van Keilegom, I. and González Manteiga, W. (2004). Comparison of regression curves based on the estimation of the error distribution. IAP-statistics Technical Report Series TR # 0417.
- Simar, L. and Zelenyuk, V. (2004). On Testing Equality of Distributions of Technical Efficiency Scores. IAP-statistics Technical Report Series TR # 0449.
- Vandenhende, F. and Lambert, P. (2004). Local dependence estimation using non-parametric archimedean copulas. IAP-statistics Technical Report Series TR # 0402.
- Van Keilegom, I. and Carroll, R. (2005). Backfitting versus profiling in general criterion functions. IAP-statistics Technical Report Series TR # 0508.
- Van Keilegom, I., González Manteiga, W. and Sánchez Sellero, C. (2004). Goodness-of-fit tests in parametric regression based on the estimation of the error distribution. IAP-statistics Technical Report Series TR # 0409.
- Van Keilegom, I. and Veraverbeke, N. (2005). U-quantile estimation in the presence of auxiliary information (in preparation).
- Wang, L. and Van Keilegom, I. (2004). Nonparametric test for the form of parametric regression with time series errors. IAP-statistics Technical Report Series TR # 0414.
- Zelenyuk, V. and Zheka, V. (2004). Corporate governance and firm’s efficiency: the case of a transitional country, Ukraine. IAP-statistics Technical Report Series TR # 0432.

## **B. List of Publications**

- Aerts, M., Claeskens, G. and Hart, J. (2004). Bayesian-motivated tests of function fit and their asymptotic frequentist properties. *Annals of Statistics*, 32, 2580-2615.
- Bauwens, L., Giot, P., Grammig, J. and Veredas, D. (2004). A comparison of financial duration models via density forecast, *International Journal of Forecasting*, 20, 589-604.

- Bauwens, L. and Veredas, D. (2004). The stochastic conditional duration model: a latent factor model for the analysis of financial durations, *Journal of Econometrics*, 119, 381-412.
- Beguín, Cl. and Simar, L. (2004). Analysis of the Expenses Linked to Hospital Stays: How to Detect Outliers, *Health Care Management Science*, 7, 89–96.
- Bigot J. (2005). A scale-space approach with wavelets to singularity estimation, accepted in *ESAIM:PES*.
- Bouezmarni, T. and Rolin, J.-M. (2004). Consistency of the beta kernel density function estimator. *The Canadian Journal of Statistics*, 31, 89-98.
- Brouhns, N., Denuit, M. and Van Keilegom, I. (2004). Bootstrapping the Poisson log-bilinear model for mortality forecasting. *Scand. Actuar. J.* (to appear).
- Claeskens, G. (2004). Restricted likelihood ratio lack of fit tests using mixed spline model. *JRSS-B*, 66, 909-926.
- Claeskens, G. and Hjort, N.L. (2004). Goodness of fit via nonparametric likelihood ratios. *Scand. J. Stat*, 31, 487-513.
- Daouia, A. and Simar, L. (2005). Robust Nonparametric Estimators of Monotone Boundaries, *Journal of Multivariate Analysis* (to appear).
- Daraio, C. and Simar, L. (2005a). Introducing environmental variables in nonparametric frontier models: a probabilistic approach, *Journal of Productivity Analysis* (to appear).
- Delaigle, A. and Gijbels, I. (2004a). Practical bandwidth selection in deconvolution kernel density estimation. *Computational Statistics and Data Analysis*, 45, 249–267.
- Delaigle, A. and Gijbels, I. (2004b). Bootstrap bandwidth selection in kernel density estimation from a contaminated sample. *The Annals of the Institute of Statistical Mathematics*, 56, 19–47.
- Delouille, V. , Simoens, J. and von Sachs, R. (2004). Smooth design-adapted wavelets for non-parametric stochastic regression, *J. Amer. Stat. Assoc*, 99, 643-658.
- Delouille, V. and von Sachs, R. (2004). Estimation of nonlinear autoregressive models using design-adapted wavelets, *Ann. Inst. Statistical Mathematics* (to appear).
- De Macq, I. and Simar, L. (2005) Hyperrectangular Space Partitioning Trees, a practical approach, *Computational Statistics* (to appear).
- Denuit, M., and Lang, S. (2004). Nonlife ratemaking with Bayesian GAM's. *Insurance: Mathematics and Economics*, 35, 627-647.
- Florens, J.P. and Simar, L. (2005). Parametric Approximations of Nonparametric Frontier. *Journal of Econometrics*, 124, 1, 91-116.



- Gijbels, I. (2004). Monotone regression. *Encyclopedia of Statistical Sciences*. Editors, S. Kotz, N.L. Johnson, C.B. Read, N. Balakrishnan and B. Vidakovic. Wiley, New York (to appear).
- Gijbels, I. and Goderniaux, A.-C. (2004a). Bandwidth selection for change point estimation in nonparametric regression. *Technometrics*, 46, 76-86.
- Gijbels, I. and Goderniaux, A.-C. (2004b). Data-driven discontinuity detection in derivatives of a regression function. *Communications in Statistics–Theory and Methods*, 33, 4, 851-871.
- Gijbels, I. and Goderniaux, A.-C. (2004c). Bootstrap test for change points in nonparametric regression. *Journal of Nonparametric Statistics*, 16, 591-611.
- Gijbels, I., Hall, P. and Kneip, A. (2004). Interval and band estimation for curves with jumps. *Journal of Applied Probability*, Special Volume 41A “Stochastic Methods and Their Applications”, Papers in honour of Chris Heyde, Edited by J. Gani and E. Seneta, 65-79.
- Gijbels, I. and Heckman, N. (2004). Nonparametric testing for a monotone hazard function via normalized spacings. *Journal of Nonparametric Statistics*, 16, 463-477.
- Hall, P. and Van Keilegom, I. (2004). Testing for monotone increasing hazard rate. *Ann. Statist.* (to appear).
- Hoeffelman J., Decat, G., Lilien, JL., Delaigle, A. and Govaerts, B. (2004). Assessment of the Electric and Magnetic field levels in the vicinity of the HV overhead power lines in Belgium, Report of the 2004 session of the International council on large electric systems, Paris.
- Lambert, P., Collett, D., Kimber, A. and Johnson, R. (2004). Parametric accelerated failure time models with random effects and an application to kidney transplant survival. *Statistics in Medicine*, 23, 3177-3192.
- Mouchart, M. and Rombouts, J. (2005). On clustered panel data models: an efficient approach for nowcasting from poor data. *International Journal of Forecasting* (to appear).
- Pitrebois, S., Denuit, M., and Walhin, J.-F. (2004). Bonus-malus scales in segmented tariffs: Gilde and Sundt’s work revisited. *Australian Actuarial Journal*, 10, 107-125.
- Purcaru, O., Guillen, M., and Denuit, M. (2004). Linear credibility models based on time series for claim counts. *Belgian Actuarial Bulletin*, 4, 62-74.
- Sánchez Sellero, C., González Manteiga, W. and Van Keilegom, I. (2004). Uniform representation of product-limit integrals with applications. *Scand. J. Statist.* (to appear).
- Simar, L. and Wilson, P.W. (2004). Performance of the bootstrap for DEA estimators and iterating the principle . International Series in Operations Research & Management Sciences: Handbook on Data Envelopment Analysis edited by W.W. Cooper, L.M. Seiford and J. Zhu, 10, 265-298.
- Van Bellegem, S. and von Sachs, R. (2004a). Forecasting economic time series with unconditional time-varying variance . *International Journal of Forecasting*, 20, 611-627.

Van Belleghem, S. and von Sachs, R. (2004b). On adaptive estimation for locally stationary wavelet processes and its applications. *International Journal of Wavelets, Multiresolution and Information Processing*, 2, 4, 545-565.

Van Keilegom, I. (2004). A note on the nonparametric estimation of the bivariate distribution under dependent censoring. *J. Nonpar. Statist.*, 16, 659-670.

### 3.1.2 Katholieke Universiteit Leuven, KUL-1 partner

#### A. List of Technical Reports

Balazs, K., Hidegkuti, I. and De Boeck, P. (2004). Detecting heterogeneity in logistic regression models. IAP-statistics Technical Report Series TR # 0470.

Ceulemans, E. and Van Mechelen, I. (2004a). An evaluation of various initial configurations and simulated annealing-based algorithmic variants for hierarchical classes analysis. IAP-statistics Technical Report Series TR # 0472.

De Fraine, B., Van Damme, J. and Onghena, P. (2004). Predicting longitudinal trajectories of adolescent academic self-concept: An application of growth mixture models. IAP-statistics Technical Report Series TR # 0474.

Depril, D. and Van Mechelen, I. (2004). One-mode additive clustering of multiway data. IAP-statistics Technical Report Series TR # 0475.

Leenen, I., Van Mechelen, I., Gelman, A. and De Knop, S. (2004). Bayesian hierarchical classes analysis. IAP-statistics Technical Report Series TR # 0476.

Prinzie, P., Onghena, P. and Hellinckx, W. (2004). A multivariate latent growth curve analysis of children's aggressive and delinquent problem behavior: Associations with harsh discipline and gender. IAP-statistics Technical Report Series TR # 0477.

San Martin, E., Del Pino, G. and De Boeck, P. (2004). IRT models for ability-based guessing. IAP-statistics Technical Report Series TR # 0479.

Schepers, J. and Van Mechelen, I. (2004). Three-mode partitioning: An evaluation of three algorithms. IAP-statistics Technical Report Series TR # 0480.

#### B. List of Publications

Ceulemans, E. and Van Mechelen, I. (2004b). Hierarchical classes models for three-way three-mode binary data: Interrelations and model selection. *Psychometrika* (to appear).

Ceulemans, E. and Van Mechelen, I. (2004c). Tucker2 hierarchical classes analysis. *Psychometrika*, 69, 375-399.

Ceulemans, E., Van Mechelen, I. and Kuppens, P. (2004). Adapting the formal to the substantive: Constrained Tucker3-HICLAS. *Journal of Classification*, 21, 19-50.

- De Boeck, P. and Smits, D. (2004). A double-structure structural equation model for the study of emotions and their components (keynote lecture). Proceedings of the XXVIIIth International Congress of Psychology (Beijing, August 2004). Psychology Press (to appear).
- De Boeck, P. and Wilson, M. (Eds.) (2004). Explanatory item response model: A generalized linear and nonlinear approach. N.Y.: Springer.
- De Boeck, P., Wilson, M. and Acton, S. (2005). A conceptual and psychometric framework for distinguishing categories and dimensions. *Psychological Review*, 112, 129-158.
- Gelman, A., Van Mechelen, I., Verbeke, G., Heitjan, D.F. and Meulders, M. (2004). Multiple imputation for model checking: Completed-data plots with missing and latent data. *Biometrics* (to appear).
- Ip, E.H., Wang, Y.J., De Boeck, P. and Meulders, M. (2004). Locally dependent latent trait models for polytomous responses. *Psychometrika*, 69, 191-216.
- Janssen, R., Schepers, J. and Peres, D. (2004). Models with item and item group predictors. In *Explanatory item response models: A generalized linear and nonlinear approach*. De Boeck, P. and Wilson, M. (Eds). New York, Springer, 189-212.
- Leenen, I. and Van Mechelen, I. (2004). A conjunctive parallelogram model for pick any/n data. *Psychometrika*, 69, 401-420.
- Meulders, M., De Boeck, P., Van Mechelen, I. and Gelman (2004). Probabilistic feature analysis of facial perception of emotions. *Applied Statistics*. (to appear).
- Meulders, M., Ip, E. and De Boeck, P. (2004). Latent variable models for partially ordered responses and trajectory analysis of anger-related feelings. *British Journal of Mathematical and Statistical Psychology* (to appear).
- Meulders, M. and Xie, Y. (2004). Person-by-item predictors. In *Explanatory item response models: A generalized linear and nonlinear approach*. De Boeck, P. and Wilson, M.(Eds). New York, Springer, 213-240.
- Rijmen, F. and Briggs, D. (2004). Multidimensional person variance and latent item predictors. In *Explanatory item response models; A generalized linear and nonlinear approach*. De Boeck, P. and Wilson, M. (eds.). New York, Springer, 247-265.
- Rijmen, F. and De Boeck, P. (2004). A relation between a between-item multidimensional IRT model and the mixture-Rasch model. *Psychometrika* (to appear).
- Rijmen, F., De Boeck, P. and van der Maas, H.L.J. (2004). An IRT model with a parameter-driven process for change. *Psychometrika*, (to appear).
- Smits, D. and De Boeck, P. (2004). The inhibition of verbally aggressive behavior. *European Journal of Personality*, 18, 537-555.

- Smits, D. J. M. and Moore, S. (2004). Latent item predictors with fixed effects. In *Explanatory item response models: A generalized linear and nonlinear approach*. De Boeck, P. and Wilson, M. (eds.). New York, Springer, 267-287.
- Tuerlinckx, F. (2004). A multivariate counting process with Weibull distributed first-arrival times. *Journal of Mathematical Psychology*, 48, 65-79.
- Tuerlinckx, F. (2004). The efficient computation of the distribution function of the diffusion process. *Behavior Research Methods, Instruments, & Computers*. 36, 702-716.
- Tuerlinckx, F. and De Boeck, P. (2004). Two interpretations of the discrimination parameter. *Psychometrika* (to appear).
- Tuerlinckx, F. and De Boeck, P. (2004). Models for residual dependencies. In , *Explanatory item response models: A generalized linear and nonlinear approach*. De Boeck, P. and Wilson, M. (Eds). New York, Springer, 289-316.
- Tuerlinckx, F., Rijmen, F., Molenberghs, G., Verbeke, G., Briggs, D., Van den Noortgate, W., Meulders, M. and De Boeck, P. (2004). Estimation and software. In: *Explanatory Item Response Models. A Generalized Linear and Nonlinear Approach*. De Boeck, P. and Wilson, M. (Eds). New York: Springer-Verlag, 343-373.
- Tuerlinckx, F. and Wang, W. C. (2004). Models for residual dependencies. In *Explanatory item response models: A generalized linear and nonlinear approach*. De Boeck, P. and Wilson, M. (eds.). New York, Springer, 289-316.
- Van den Noortgate, W. and De Boeck, P. (2004). Assessing and explaining differential item functioning (DIF) using logistic mixed models. *Journal of Educational and Behavioral statistics*, to appear.
- Van Den Noortgate, W. and Paek, I. (2004). Person regression models. . In *Explanatory item response models: A generalized linear and nonlinear approach*. De Boeck, P. and Wilson, M.(eds.). New York, Springer, 167-187.
- Van Mechelen, I., Bock, H.-H. and De Boeck, P. (2004a). Two-mode clustering methods: A structured overview. *Statistical Methods in Medical Research*, 12, 363-394.
- Van Mechelen, I., Bock, H.-H. and De Boeck, P. (2004b). Two-mode clustering methods. In: Everitt, David Howell (Eds.), *Encyclopedia Statistics in Behavioral Science*. N.Y., Wiley.
- Wilson, M. and De Boeck, P. (2004). Descriptive and explanatory item response models. In: *Explanatory item response models: A generalized linear and nonlinear approach*. De Boeck, P. and Wilson, M. (eds.). New York: Springer, 43-74.

### 3.1.3 Katholieke Universiteit Leuven, KUL-2 partner

#### A. List of Technical Reports

- Fieuws, S. and Verbeke, G. (2004). Pairwise fitting of mixed models for the joint modelling of multivariate longitudinal profiles. IAP-statistics Technical Report Series TR # 0450.
- Komárek, A. and Lesaffre, E. (2004a). A Bayesian accelerated failure time model with a normal mixture as an error distribution. IAP-statistics Technical Report Series TR # 0452.
- Komárek, A. and Lesaffre, E. (2004b). Bayesian accelerated failure time model for correlated censored data with a normal mixture as an error distribution. IAP-statistics Technical Report Series TR # 0482.
- Lesaffre, E. and Bogaerts, K. (2004). Estimating Kendall's tau for bivariate interval censored data with a smooth estimate of the density. IAP-statistics Technical Report Series TR # 0454.
- Lesaffre, E., Rizopoulos, D. and Tsonaka, S. (2004). The logistic-transform for bounded outcome scores. IAP-statistics Technical Report Series TR # 0448.
- Mwalili, S.M., Lesaffre, E. and Declerck, D. (2004). The interval censored zero-inflated negative binomial regression model: an application in caries research. IAP-statistics Technical Report Series TR # 0462.
- Tsonaka, S. Rizopoulos, D. and Lesaffre, E. (2005). Power and sample size calculations for discrete bounded outcome scores. IAP-statistics Technical Report Series TR # 0505.

#### B. List of Publications

- Adair, P.M., Pine, C.M., Burnside, G., Nicoll, A.D., Gillett, A., Anwar, S., Broukal, Z., Chestnutt, I.G., Declerck, D., Ping, F.X., Ferro, R., Freeman, R., Grant-Mills, D., Gugushe, T., Hunsrisakhun, J., Irigoyen-Camacho, M., Lo, E.C., Moola, M.H., Naidoo, S., Nyandindi, U., Poulsen, V.J., Ramos-Gomez, F., Razanamihaja, N., Shahid, S., Skeie, M.S., Skur, O.P., Splieth, C., Soo, T.C., Whelton, H. and Young, D.W. (2004). Familial and cultural perceptions and beliefs of oral hygiene and dietary practices among ethnically and socio-economically diverse groups. *Community Dental Health*, 21(Supplement), 102-111.
- Bogaerts, K., Leroy R., Lesaffre E. and Declerck, D. (2004). Modelling tooth emergence data based on multivariate interval-censored data. *Stat Med.* 30, 3775-3787.
- Bogaerts, K. and Lesaffre, E. (2004a). A smooth estimate of the bivariate survival density in the presence of left, right and interval censored data. In Proceedings of the Joint Statistical Meetings, Biometrics Section, 633-639, Alexandria, VA: American Statistical Association.
- Bogaerts, K. and Lesaffre, E. (2004b). A new, fast algorithm to find the regions of possible support for bivariate interval-censored data. *Journal of Computational and Graphical Statistics*, 13, 330-340.

- Bottenberg, P., Declerck, D., Ghidey, W., Bogaerts, K., Vanobbergen, J. and Martens, L. (2004). Prevalence and determinants of enamel fluorosis in Flemish schoolchildren. *Caries Research*, 38, 20-28.
- Gelman, A., Van Mechelen, I., Verbeke, G., Heitjan, D.F. and Meulders, M. (2004). Multiple imputation for model checking: Completed-data plots with missing and latent data. *Biometrics* (to appear).
- Ghidey, W., Lesaffre, E. and Eilers, P. (2004). Smooth random effects distribution in a linear mixed model. *Biometrics*, 60, 945-953.
- Hens, N., Aerts, M., Molenberghs, G., Thijs, H. and Verbeke, G. (2004). Kernel weighted influence measures. *Computational Statistics and Data Analysis* (to appear).
- Komárek, A., Lesaffre, E., Härkänen, T., Declerck, D. and Virtanen, J.I. (2005). A Bayesian analysis of multivariate doubly interval censored dental data. *Biostatistics*, 6, 145-155.
- Komárek, A., Lesaffre, E. and Hilton, J.F. (2005). Accelerated failure time model for arbitrarily censored data with smoothed error distribution. *Journal of Computational and Graphical Statistics* (to appear).
- Lesaffre, E., Komárek, A. and Declerck, D. (2005). An overview of methods for interval-censored data with an emphasis on applications in dentistry. *Statistical Methods in Medical Research* (to appear).
- Lesaffre, E., Kocmanová, D., Lemos, P.A. Disco, C.M.C. and Serruys, P.W. (2003). A retrospective analysis of the effect of noncompliance on time to first major adverse cardiac events in LIPS. *Clinical Therapeutics*, 25, 2431-2447.
- Lesaffre, E., Mwalili, S.M. and Declerck, D. (2004). Analysis of caries experience taking inter-observer bias and variability into account. *Journal of Dental Research*, 83, 951-955.
- Martens, L., Vanobbergen, J., Leroy, R., Lesaffre, E. and Declerck, D. (2004). Variables associated with oral hygiene levels in 7-year-olds in Belgium. *Community Dental Health*, 21, 4-10.
- Molenberghs, G., Thijs, H., Michiels, B., Verbeke, G. and Kenward, M.G. (2004). Pattern-mixture models. *Journal de la Société française de Statistique*, 145, 49-77.
- Molenberghs, G. and Verbeke, G. (2004a). Meaningful statistical model formulations. *Statistica Sinica*, 14, 177-206.
- Molenberghs, G. and Verbeke, G. (2004b). An introduction to (generalized) (non-)linear mixed models. In: *Explanatory Item Response Models. A Generalized Linear and Nonlinear Approach*. De Boeck, P. and Wilson, M. (Eds). New York: Springer-Verlag, 111-153.
- Mwalili, S.M., Lesaffre, E. and Declerck, D. (2005). A Bayesian ordinal logistic regression model to correct for interobserver measurement error in a geographical oral health study. *Applied Statistics*, 54, 77-93.

Tuerlinckx, F., Rijmen, F., Molenberghs, G., Verbeke, G., Briggs, D., Van den Noortgate, W., Meulders, M. and De Boeck, P. (2004). Estimation and software. In: *Explanatory Item Response Models. A Generalized Linear and Nonlinear Approach*. De Boeck, P. and Wilson, M. (Eds). New York: Springer-Verlag, 343-373.

Vanobbergen, J., Declerck, D., Mwalili, S. and Martens, L. (2004). The effectiveness of a 6-year oral health education programme for primary schoolchildren. *Community Dentistry and Oral Epidemiology*, 32, 173-82.

### 3.1.4 Limburgs Universitair Centrum, LUC partner

#### A. List of Technical Reports

Cadarso-Suarez, C., Roca-Pardinas, J., Molenberghs, G., Faes, C., Nacher, V., Ojeda, S. and Acuna, C. (2004). Flexible modeling of neuron firing rates across different experimental conditions. An application to neural activity in the prefrontal cortex during a discrimination task. IAP-statistics Technical Report Series TR # 0463.

Cortiñas, J., Legrand, C., Burzykowski, T., Duchateau, L. and Janssen, P. (2004). Comparison of different estimation procedures used to deal with proportional hazards model with random effects. IAP-statistics Technical Report Series TR # 0407.

Cortiñas Abrahantes, J. and Burzykowski, T. (2004). A version of the EM algorithm for proportional hazard model with random effects. IAP-statistics Technical Report Series TR # 0455.

Cortiñas Abrahantes, J., Legrand, Burzykowski, T., Janssen, Ducrocq, and Duchateau, L. (2004). Comparison of different estimation procedures for proportional hazards model with random effect. IAP-statistics Technical Report Series TR # 0456.

Duchateau, L. and Janssen, P. (2004). Understanding heterogeneity in mixed, generalized mixed and frailty models. IAP-statistics Technical Report Series TR # 0412.

Faes, C., Geys, H., Aerts, M. and Molenberghs, G. (2004). A hierarchical modeling approach for risk assessment in developmental toxicity studies. IAP-statistics Technical Report Series TR # 0464.

Faes, C., Geys, H., Molenberghs, G., Aerts, M., Cadarso-Suarez, C., Acuna, C. and Cano, M. (2004). A flexible method to measure synchrony in neuronal firing. IAP-statistics Technical Report Series TR # 0467.

Hens, N., Aerts, M. and Molenberghs, G. (2004). Model selection for incomplete and design-based samples. IAP-statistics Technical Report Series TR # 0466.

Janssen, P., Swanepoel, J. and Veraverbeke, N. (2004). Modifying the kernel distribution function estimator towards reduced bias. IAP-statistics Technical Report Series TR # 0458.

- Legrand, C., Duchateau, L., Sylvester, R., Janssen, P., van der Hage, J., van de Velde, C., and Therasse, P. (2004). Heterogeneity in outcome between institutions: lessons learned from an EORTC cancer trial. Revision of IAP-statistics Technical Report Series TR # 0405.
- Legrand, Ducrocq, Janssen, Sylvester and Duchateau, L. (2004). A bayesian approach to jointly estimate center and treatment by center heterogeneity in a proportional hazards model. IAP-statistics Technical Report Series TR # 0457.
- Massonnet, G., Burzykowksi, T. and Janssen, P. (2004). Resampling plans for frailty models. Revision of IAP-statistics Technical Report Series TR # 0348.
- Nguti, R., Claeskens, G. and Janssen, P. (2004). One-sided tests in shared frailty models. IAP-statistics Technical Report Series TR # 0436.
- Van Keilegom, I. and Veraverbeke, N. (2005). U-quantile estimation in the presence of auxiliary information (in preparation).
- Veraverbeke, N. (2004). Conditional residual lifetime under random censorship. IAP-statistics Technical Report Series TR # 0406.
- Wang, L. (2004a). A note on the choice between two loss functions in Bayesian analysis. IAP-statistics Technical Report Series TR # 0438.
- Wang, L. (2004b). Bayes and empirical bayes tests for the life parameter. IAP-statistics Technical Report Series TR # 0440.
- Wang, L. and Veraverbeke, N. (2004a). Bayes prediction under random censorship. IAP-statistics Technical Report Series TR # 0439.
- Wang, L. and Veraverbeke, N. (2004b). Empirical likelihood in a semiparametric model for missing response data. IAP-statistics Technical Report Series TR # 0441.

## B. List of Publications

- Aerts, C, Lamers, H.J.G.L.M. and Molenberghs, G. (2004). Maximum mass-loss rates of line-driven winds of massive stars: the effect of rotation and an application to h Carinae. *Astronomy and Astrophysics*, 418, 639-648.
- Aerts, M., Claeskens, G. and Hart, J. (2004). Bayesian-motivated tests of function fit and their asymptotic frequentist properties. *Annals of Statistics*, 32, 2580-2615.
- Alonso, A., Geys, H., Molenberghs, G., Kenward, M.G. and Vangeneugden, T. (2004). Validation of surrogate markers in multiple randomized clinical trials with repeated measurements: canonical correlation approach. *Biometrics*, 60, 845-853.
- Alonso, A., Molenberghs, G., Burzykowski, T., Renard, D., Geys, H., Shkedy, Z., Tibaldi, F., Cortiñas, J. and Buyse, M. (2004). Prentice's approach and the meta-analytic paradigm: a reflection on the role of statistics in the evaluation of surrogate endpoints. *Biometrics*, 60, 724-728.



- Burzykowski, T., Molenberghs, G. and Buyse, M. (2004). The validation of surrogate endpoints using data from randomized clinical trials: a case-study in advanced colorectal cancer. *Journal of the Royal Statistical Society, Series A*, 167, 103-124.
- Cao, R., Janssen, P. and Veraverbeke, N. (2004). Relative hazard rate estimation for censored and left truncated data. *TEST*.
- Cao, R., Lopez-De-Ullibarri, I., Janssen, P. and Veraverbeke, N.. (2004). Efficiency of pre-smoothed Kaplan-Meier and Nelson-Aalen estimators . *Nonpar. Statist.*
- Cortiñas Abrahantes, J., Molenberghs, G., Burzykowski, T., Shkedy, Z. and Renard, D. (2004). Choice of units of analysis and modeling strategies in multilevel hierarchical models. *Computational Statistics and Data Analysis*, 47, 537-563.
- Cserni, G., Burzykowski, T., Vinh-Hung, V., Kocsis, L., Boross, G., Sinko, M., Tarjan, M., Bori, R., Rajtar, M., Tekle, E., Maraz, R., Baltas , B. and Svebis, M. (2004). Axillary sentinel node and tumour-related factors associated with non-sentinel node involvement in breast cancer. *Japanese Journal of Clinical Oncology*, 34, 519-524.
- Decin, L., Shkedy, Z., Molenberghs, G., Aerts, C. and Aerts, M. (2004). Estimating stellar parameters from spectra: I. Non-parametric estimator for the spectrum and lack-of-fit test. *Astronomy and Astrophysics*, 421, 281-294.
- De Ketelaere, B., Lammertyn, J., Molenberghs, G., Desmet, M., Nicola, B. and De Baerdemaeker, J.(2004). Tomato cultivar grouping based on firmness evolution, shelf-life and variability during postharvest storage. *Postharvest Biology and Technology*, 34, 187-201.
- Dmitrienko, A., Offen, W.W., Faries, D., Christy Chuang-Stein, J.L., and Molenberghs, G. (2004). *Analysis of Clinical Trial Data Using the SAS System*. Cary, NC: Sas Publishing.
- Duchateau, L. and Janssen, P. (2004). Penalized partial likelihood for frailties and smoothing splines in time to first insemination models in dairy cows. *Biometrics*, 60, 608-614.
- Duchateau, L., Opsomer, G., Dewulf, J. and Janssen, P. (2005). The nonlinear effect (determined by the penalised partial likelihood approach) of milk-protein concentration on time to first insemination in Belgian diary cows. *Preventive Veterinary Medicine* (to appear).
- Faes, C., Aerts, M., Geys, H., Molenberghs, G. and Declerck, L. (2004). Bayesian testing for trend in a power model for clustered binary data. *Environmental and Ecological Statistics*, 11, 305-322.
- Faes, C., Geys, H., Aerts, M., Molenberghs, G. and Catalano, P. (2004). Modeling combined continuous and ordinal outcomes in a clustered setting. *Journal of Agricultural, Biological, and Environmental Statistics*, 9, 1-16.
- Hens, N., Aerts, M., Molenberghs, G., Thijs, H. and Verbeke, G. (2004). Kernel weighted influence measures. *Computational Statistics and Data Analysis* (to appear).

- Katsambas, A., Abeck, D., Haneke, E., van de Kerkhof, P., Burzykowski, T., Molenberghs, G. and Marynissen, G. (2004). The effects of foot disease on quality of life: results from the Achilles Project. *Journal of the European Academy of Dermatology and Venerology* (to appear).
- Lammers, W., Faes, C., Stephen, B., Bijmens, L., Ver Donck, L. and Schuurkes, J.A.J. (2004). Spatial determination of successive spikes in the isolated cat duodenum. *Neurogastroenterology and motility*, 16, 775-783.
- Lipsitz, S.R., Molenberghs, G., Fitzmaurice, G.M. and Ibrahim, J.G. (2004). A protective estimator for linear regression with nonignorably missing Gaussian outcomes. *Statistical Modelling*, 4, 3-17.
- Mallinckrodt, C.H., Watkin, J.G., Molenberghs, G. and Carroll, R.J. (2004). Choice of the primary analysis in longitudinal clinical trials. *Pharmaceutical Statistics*, 3, 161-169.
- Mallinckrodt, C.H., Kaiser, C.J., Watkin, J.G., Molenberghs, G. and Carroll, R.J. (2004). Type I error rates from likelihood-based repeated measures analysis of incomplete longitudinal data. *Pharmaceutical Statistics*, 3, 171-186.
- Molenberghs, G., Burzykowski, T., Alonso, A. and Buyse, M. (2004). A perspective on surrogate endpoints in controlled clinical trials. *Statistical Methods in Medical Research*, 13, 177-206.
- Molenberghs, G., Buyse, M. and Burzykowski, T. (2004). Surrogate Markers. *Bulletin of the International Chinese Statistical Association* (to appear).
- Molenberghs, G., Thijs, H., Jansen, I., Beunckens, C., Kenward, M.G., Mallinckrodt, C. and Carroll, R.J. (2004). Analyzing incomplete longitudinal clinical trial data. *Biostatistics*, 5, 445-464.
- Molenberghs, G., Thijs, H., Michiels, B., Verbeke, G. and Kenward, M.G. (2004). Pattern-mixture models. *Journal de la Société française de Statistique*, 145, 49-77.
- Molenberghs, G. and Verbeke, G. (2004a). Meaningful statistical model formulations. *Statistica Sinica*, 14, 177-206.
- Molenberghs, G. and Verbeke, G. (2004b). An introduction to (generalized) (non-)linear mixed models. In: *Explanatory Item Response Models. A Generalized Linear and Nonlinear Approach*. De Boeck, P. and Wilson, M. (Eds). New York: Springer-Verlag, 111-153.
- Moons, E., Aerts, M. and Wets, G. (2004). A tree based lack-of-fit test for multiple logistic regression. *Statistics in Medicine*, 23, 1425-1438.
- Renard, D., Molenberghs, G. and Geys, H. (2004). A pairwise likelihood approach to estimation in multilevel probit models. *Computational Statistics and Data Analysis*, 44, 649-667.
- Sakamoto, J., Ohashi, Y., Hamada, C., Buyse, M., Burzykowski, T. and Piedbois, P. (2004). Efficacy of oral adjuvant therapy after resection of colorectal cancer: 5-year results from three randomized trials. *Journal of Clinical Oncology*, 22, 484-492.

- Schrooten, W., Florence, E., Dreezen, C., Van Esbroeck, M., Fransen, K., Alonso, A., Desmet, P., Colebunders, R., Kestens, L. and Roo De, A. (2004). Five-year immunological outcome of highly antiretroviral treatment in a clinical setting: results from a single HIV treatment center. *International Journal of STD and AIDS*, 15, 523-528.
- Shkedy, Z., Molenberghs, G., Van Craenendonck, H., Aerts, N., Steckler, T. and Bijmens, L. (2004). A hierarchical binomial-Poisson model for the analysis of a cross-over design for correlated binary data when the number of trials is dose-dependent. *Journal of Biopharmaceutical Statistics* (to appear).
- Speybroeck, N., Berkvens, D., Mfoukou-Ntsakala, A., Aerts, M., Hens, N., Van Huylenbroeck, G. and Thys, E. (2004). Classification trees versus multinomial models in the analysis of urban farming systems in Central Africa. *Agricultural Systems*, 80, 133-149.
- Tibaldi, F., Molenberghs, G., Burzykowski, T. and Geys, H. (2004). Pseudo-likelihood estimation for a marginal multivariate survival model. *Statistics in Medicine*, 23, 947-963.
- Tibaldi, F., Torres Barbosa, F. and Molenberghs, G. (2004). Modelling associations between time-to-event responses in pilot cancer clinical trials using a Plackett-Dale model. *Statistics in Medicine*, 23, 2173-2186.
- Tuerlinckx, F., Rijmen, F., Molenberghs, G., Verbeke, G., Briggs, D., Van den Noortgate, W., Meulders, M. and De Boeck, P. (2004). Estimation and software. In: *Explanatory Item Response Models. A Generalized Linear and Nonlinear Approach*. De Boeck, P. and Wilson, M. (Eds). New York: Springer-Verlag, 343-373.
- Vangeneugden, T., Laenen, A., Geys, H., Renard, D. and Molenberghs, G. (2004). Applying linear mixed models to estimate reliability in clinical trial data with repeated measurements. *Controlled Clinical Trials*, 25, 13-30.
- Van Steen, K., Tahri, N. and Molenberghs, G. (2004). Introducing the multivariate Dale model in population-based genetic association studies. *Biometrical Journal*, 46, 187-202.

### **3.1.5 Université Libre de Bruxelles, ULB partner**

#### **A. List of Technical Reports**

- Ayadi, A. and Mélard, G. (2004). On the distribution of the sample autocorrelations of a white noise process. IAP-statistics Technical Report Series TR # 0461.
- Bramati, M.C. and Croux, C. (2005). Robust estimators for the fixed effects panel data model. IAP-statistics Technical Report Series TR # 0512.
- Cohen, A., Mélard, G. and Ouakasse, A. (2004). Une expérience de télé-enseignement en statistique pour une banque centrale: aspects technologiques. IAP-statistics Technical Report Series TR # 0446.
- Croux, C. and Dehon, C. (2005). Robustness versus efficiency for nonparametric correlation measures. IAP-statistics Technical Report Series TR # 0513.

- D'Agostino, A. and Giannone, D. (2005). Comparing alternative predictors based on large-panel dynamic factor models. IAP-statistics Technical Report Series TR # 0514.
- Defrise, M. and De Mol, C. (2005). Linear inverse problems with mixed smoothness and sparsity constraints. IAP-statistics Technical Report Series TR # 0515.
- Dehon, C. and Fiszman, P. (2005). Social and spatial inequalities in consulting a GP versus a specialist in Belgium. IAP-statistics Technical Report Series TR # 0516.
- Dolado, J., Rodriguez Poo, L. and Veredas, D. (2004). Testing weak exogeneity in the exponential family: an application to financial point processes. IAP-statistics Technical Report Series TR # 0483.
- Forni, M., Giannone, D., Lippi, M. and Reichlin, L. (2004). Opening the Black Box: Structural factor models vs structural VARs, mimeo. IAP-statistics Technical Report Series TR # 0484.
- Giannone, D. and Lenza, M. (2004). The Feldstein-Horioka fact. IAP-statistics Technical Report Series TR # 0485.
- Hallin, M., Jureckova, J. and Koul, H. L. (2005). Serial Autoregression and Regression Rank Scores Statistics. IAP-statistics Technical Report Series TR # 0507.
- Hallin, M., Oja, H. and Paindaveine, D. (2004). Optimal R-estimation of shape. IAP-statistics Technical Report Series TR # 0430.
- Hallin, M. and Saidi, A. (2004b). Optimal tests for non-correlation between multivariate time series. IAP-statistics Technical Report Series TR # 0428.
- Hallin, M., Vermandele, C. and Werker, B. (2004). Semiparametrically efficient inference based on signs and ranks for median restricted models. IAP-statistics Technical Report Series TR # 0403.
- Hallin, M., Vermandele, C. and Werker, B. (2003). Serial and nonserial sign-and-rank statistics. Asymptotic representation and asymptotic normality. IAP-statistics Technical Report Series TR # 0305.
- Klein, A., Mélard, G. and Niemczyk, J. (2004). Corrections to "Construction of the exact Fisher information matrix of Gaussian time series models by means of matrix differential rules". IAP-statistics Technical Report Series TR # 0460.
- Moulin, L., Salto, M., Silvestrini, A. and Veredas, D. (2004). Using intra annual information to forecast the annual state deficits. The case of France. IAP-statistics Technical Report Series TR # 0486.
- Paindaveine, D. (2004). Chernoff-Savage and Hodges-Lehmann results for Wilks' test of multivariate independence. IAP-statistics Technical Report Series TR # 0473.
- Paindaveine, D. (2004). A Chernoff-Savage result for shape. On the non-admissibility of pseudo-Gaussian methods. IAP-statistics Technical Report Series TR # 0447.

Pascual, R. and Veredas, D. (2004). What pieces of limit order book information are informative? An empirical analysis of a pure order driven market. IAP-statistics Technical Report Series TR # 0478.

## B. List of Publications

Azrak, R. and Mélard, G. (2004). Asymptotic properties of quasi-maximum likelihood estimators for ARMA models with time-dependent coefficients. *Statistical Inference for Stochastic Processes* (to appear).

Azrak, R., Mélard, G. and Njimi, H. (2004). Forecasting in the analysis of mobile telecommunication data - Correction for outliers and replacement of missing observations. *Journal Marocain d'Automatique, d'Informatique et de Traitement du Signal*, special issue CoPSTIC'03, 1-14.

Bauwens, L., Giot, P., Grammig, J. and Veredas, D. (2004). A comparison of financial duration models via density forecast, *International Journal of Forecasting*, 20, 589-604.

Bauwens, L. and Veredas, D. (2004). The stochastic conditional duration model: a latent factor model for the analysis of financial durations, *Journal of Econometrics*, 119, 381-412.

Cristadoro, R., Forni, M., Reichlin, L. and Veronese, G. (2004). A measure of core inflation for the EURO area, *Journal of Money Credit and Banking* (to appear).

Daubechies, I., Defrise, M. and De Mol, C. (2004). An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications in Pure and Applied Mathematics*, 57, 1413-1457.

Defrise, M. and De Mol, C. (2004). Inverse imaging with mixed penalties. *Proceedings URSI EMTS 2004*, PLUS Ed., Univ. Pisa, Pisa, 798-800.

Dufour, J.-M., Farhat, A. and Hallin, M. (2004). Distribution-free bounds for serial correlation coefficients in heteroskedastic symmetric time series. *Journal of Econometrics* (to appear).

Forni, M., Hallin, M., Lippi, M. and Reichlin, L. (2004). The generalized dynamic factor model: consistency and rates. *Journal of Econometrics*, 119, 231-255.

Forni, M., Hallin, M., Lippi, M. and Reichlin, L. (2005). The generalized dynamic factor model: one-sided estimation and forecasting. *Journal of the American Statistical Association* (to appear).

Giannone, D., Reichlin, L. and Sala, L. (2005a). Monetary policy in real time in *Macroeconomic Annual 2004*, M. Gertler and K. Rogoff (Eds), MIT Press 2005 (to appear).

Giannone, D., Reichlin, L. and Sala, L. (2005b). VARs, factor models and the empirical validation of equilibrium business cycle models. *Journal of Econometrics* (to appear).

Giannone, D. and Reichlin, L. (2005). Euro area and US recessions: 1970-2003, in *The Euro Area Business Cycle*, L. Reichlin (ed.), CEPR (to appear).

- Hallin, M. (2004). Journées de Statistique: la Conférence Lucien Le Cam. *Journal de la Société française de Statistique*, 145, 3-6.
- Hallin, M. and Lotfi, S. (2004). Optimal detection of periodicities in vector autoregressive models, in P. Duchesne and B. Rémillard (Eds), *Statistical Modeling and Analysis for Complex Data Problems*, Kluwer, 2004, 49-75.
- Hallin, M., Lu, Z. and Tran, L.T. (2004a). Kernel density estimation for spatial processes: the  $L_1$  theory. *Journal of Multivariate Analysis*, 88, 61-75.
- Hallin, M., Lu, Z. and Tran, L.T. (2004b). Local linear spatial regression, *Annals of Statistics*, 32, 2469-2500.
- Hallin, M. and Paindaveine, D. (2004a). Affine-invariant aligned rank tests for the multivariate general linear model with ARMA errors, *Journal of Multivariate Analysis*, 93, 122-163.
- Hallin, M. and Paindaveine, D. (2004b). Rank-based optimal tests of the adequacy of an elliptic VARMA model, *Annals of Statistics*, 32, 2642-2678.
- Hallin, M. and Paindaveine, D. (2005a). Asymptotic linearity of serial and nonserial multivariate signed rank statistics. *Journal of Statistical Planning and Inference*, forthcoming.
- Hallin, M. and Paindaveine, D. (2005b). Multivariate signed rank tests in vector autoregressive order identification. *Statistical Science* (to appear).
- Hallin, M. and Paindaveine, D. (2005c). Optimal rank-based tests for sphericity. *Annals of Statistics*, 33 (to appear).
- Hallin, M. and Saidi, A. (2004a). Testing non-correlation and non-causality between multivariate ARMA time series. *Journal of Time Series Analysis*, 26, 83-105.
- Klein, A. and Mélard, G. (2004). An algorithm for computing the asymptotic Fisher information matrix for seasonal SISO models. *Journal of Time Series Analysis*, 25, 627-648.
- Klein, A., Mélard, G. and Spreij, P. (2005). On the resultant property of the Fisher information matrix of a vector ARMA process. *Linear Algebra and its Applications* (to appear).
- Mélard, G., Roy, R. and Saidi, A. (2005). Exact maximum likelihood estimation of structured or unit root multivariate time series models. *Computational Statistics and Data Analysis* (to appear).
- Niemczyk, J. (2004). Computing the derivatives of the autocovariances of a VARMA process, in *Proceedings in Computational Statistics*, Jaromir Antoch (Ed.), Physica-Verlag, Heidelberg, 2004, 1593-1600.
- Oja, H. and Paindaveine, D. (2005). Optimal signed-rank tests based on hyperplanes. *Journal of Statistical Planning and Inference* (to appear).
- Paindaveine, D. (2004). A unified and elementary proof of serial and nonserial, univariate and multivariate, Chernoff-Savage results. *Statistical Methodology*, 1, 81-91.

Reichlin, L. (2005). *The Euro Business Cycle: Stylized Facts and Measurement Issues*, CEPR, London 2004 (to appear).

Veredas, D. and Pascual, R. (2004). Qu componentes del libro de rdenes son informativos, *Revista Bolsa de Madrid*, 31.

### **3.1.6 Aachen Technical University, RWTH partner**

#### **A. List of Technical Reports**

De Carvalho, F., Brito, P. and Bock, H.-H. (2004). Dynamic clustering for interval data based on L2 distance. Pattern Recognition. IAP-statistics Technical Report Series TR # 0437.

#### **B. List of Publications**

Beutner, E. (2004). Risk minimization in financial markets when considering transaction costs. Dissertation, RWTH Aachen.

Bock, H.-H. (2003a). Two-way clustering for contingency tables: Maximizing a dependence Measure (pp. 143-154). In: M. Schader, W. Gaul and M. Vichi (Eds.). *Between data science and applied data analysis*. Springer-Verlag, Heidelberg. 143-154.

Bock, H.-H. (2003b). Convexity-based clustering criteria: theory, algorithms, and applications in statistics. *Statistical Methods and Applications*, 12, 293-317.

Bock, H.-H. (2004). Proximity measures. In: Brian Everitt, David Howell (Eds.), *Encyclopedia Statistics in Behavioral Science*. N.Y.: Wiley.

Chiodi, M., Mineo, A. and Bock, H.-H. (2004). *Advances in multivariate data analysis*. Heidelberg: Springer Verlag (281 pp.).

Van Mechelen, I., Bock, H.-H. and De Boeck, P. (2004a). Two-mode clustering methods: A structured overview. *Statistical Methods in Medical Research*, 12, 363-394.

Van Mechelen, I., Bock, H.-H. and De Boeck, P. (2004b). Two-mode clustering methods. In: Everitt, David Howell (Eds.), *Encyclopedia Statistics in Behavioral Science*. N.Y., Wiley.

### **3.1.7 Université Joseph Fourier, UJF-LMC-IMAG partner**

#### **A. List of Technical Reports**

Amato, U., Antoniadis, A. and Pensky, M. (2004). Wavelet kernel penalized estimation for non-equispaced design regression. IAP-statistics Technical Report Series TR # 0426.

Antoniadis, A. and Bigot, J. (2004). Poisson inverse problems. IAP-statistics Technical Report Series TR # 0425.

Antoniadis, A., Bigot, J. and Gijbels, I. (2005). Penalized wavelet monotone regression (in preparation).

- Antoniadis, A., Bigot J. and von Sachs, R. (2005). Statistical analysis and characterization of brain response images (in preparation).
- Antoniadis, A., Gijbels, I. and Nikolova, M.(2005). Penalized likelihood regression for generalized linear models with nonquadratic penalties (in preparation).
- Antoniadis, A. and Fryzlewitz, P. (2004). Parametric modelling of thresholds across scales in wavelet regression. IAP-statistics Technical Report Series TR # 0481.
- Antoniadis, A. and Sapatinas, T. (2004). Estimation and inference in functional mixed-effects models. IAP-statistics Technical Report Series TR # 0459.
- Fort, G., Lambert-Lacroix, S. and Peyre, J.(2004). Réduction de dimension dans les modèles généralisés : application à la classification de données issues de biopuces. IAP-statistics Technical Report Series TR # 0471.
- Girard, S., Iouditski, A. and Nazin, A. (2005). Linear programming problems for L1 optimal frontier estimation. IAP-statistics Technical Report Series TR # 0506.

## B. List of Publications

- Abramovich, F., Antoniadis, A. Sapatinas, T. and Vidakovic, B. (2004). Optimal testing in a fixed-effects functional analysis of variance models. *Int. J. Wavelets, Multiresolution Inf. Proc.*, 2, 323–349.
- Amato, U., Antoniadis, A. and De Feis, I. (2004). Dimension reduction in functional regression with applications, *Computational Statistics and Data Analysis* (to appear).
- Antoniadis, A., Grégoire, G. and McKeague, I. (2004), Bayesian estimation in single-index models. *Stat. Sin.*, 14, 1147–1164.
- Fort, G. and Lambert-Lacroix, S. (2004a). Ridge-partial least squares for generalized linear models with binary response. In J. Antoch, editor, 16th Symposium of IASC, COMPSTAT'04, *Proceedings in Computational Statistics*, 1019–1026, Physica Verlag/Springer.
- Fort, G. and Lambert-Lacroix, S. (2004b). Classification using partial least squares with penalized logistic regression. *Bioinformatics* (to appear).
- Gardes, L. and Girard, S. (2005). Asymptotic properties of a Pickands type estimator of the extreme value index, to appear In *Focus on probability theory*, F. Colombus, editor, Nova Science, New-York (to appear).
- Le Borgne, H., Guérin, A. and Antoniadis, A. (2004). Representation of images for classification with independent features. *Pattern Recognition Lett*, 25, 141–154.
- Mercier, G., Berthault, N., Mary, J., Peyre, J., Antoniadis, A., Comet, J-P., Cornuejol, A., Froidevaux, C. and Dutreix, M. (2004). Biological detection of low radiation doses by combining results of two microarray analysis methods. *Nucleic Acids Research*, 32, 1-7.
- Sardy, S., Antoniadis, A. and Tseng, P. (2004). Automatic smoothing with wavelets for a wide class of distributions. *J. Comput. Graphical Stat*, 13, 399-421.



## 3.2 List of Joint publications

### A. List of Joint Technical Reports

- Antoniadis, A. and Bigot, J. (2004). Poisson inverse problems, IAP-statistics Technical Report Series TR # 0425.
- Antoniadis, A., Bigot, J. and Gijbels, I. (2005). Penalized wavelet monotone regression (in preparation).
- Antoniadis, A., Bigot J. and von Sachs, R. (2005). Statistical analysis and characterization of brain response images (in preparation).
- Antoniadis, A., Gijbels, I. and Nikolova, M. (2005). Penalized likelihood regression for generalized linear models with nonquadratic penalties (in preparation).
- Nguti, R., Claeskens, G. and Janssen, P. (2004). One-sided tests in shared frailty models. IAP-statistics Technical Report Series TR # 0436.
- Van Keilegom, I. and Veraverbeke, N. (2005). U-quantile estimation in the presence of auxiliary information (in preparation).

### B. List of Joint Publications

- Aerts, M., Claeskens, G. and Hart, J. (2004). Bayesian-motivated tests of function fit and their asymptotic frequentist properties. *Annals of Statistics*, 32, 2580-2615.
- Bauwens, L., Giot, P., Grammig, J. and Veredas, D. (2004). A comparison of financial duration models via density forecast, *International Journal of Forecasting*, 20, 589-604.
- Bauwens, L. and Veredas, D. (2004). The stochastic conditional duration model: a latent factor model for the analysis of financial durations, *Journal of Econometrics*, 119, 381-412.
- Gelman, A., Van Mechelen, I., Verbeke, G., Heitjan, D.F. and Meulders, M. (2004). Multiple imputation for model checking: Completed-data plots with missing and latent data. *Biometrics* (to appear).
- Hens, N., Aerts, M., Molenberghs, G., Thijs, H. and Verbeke, G. (2004). Kernel weighted influence measures. *Computational Statistics and Data Analysis* (to appear).
- Molenberghs, G., Thijs, H., Michiels, B., Verbeke, G. and Kenward, M.G. (2004). Pattern-mixture models. *Journal de la Société française de Statistique*, 145, 49-77.
- Molenberghs, G. and Verbeke, G. (2004a). Meaningful statistical model formulations. *Statistica Sinica*, 14, 177-206.
- Molenberghs, G. and Verbeke, G. (2004b). An introduction to (generalized) (non-)linear mixed models. In: *Explanatory Item Response Models. A Generalized Linear and Nonlinear Approach*. De Boeck, P. and Wilson, M. (Eds). New York: Springer-Verlag, 111-153.

- Tuerlinckx, F., Rijmen, F., Molenberghs, G., Verbeke, G., Briggs, D., Van den Noortgate, W., Meulders, M. and De Boeck, P. (2004). Estimation and software. In: *Explanatory Item Response Models. A Generalized Linear and Nonlinear Approach*. De Boeck, P. and Wilson, M. (Eds). New York: Springer-Verlag, 343-373.
- Van Mechelen, I., Bock, H.-H. and De Boeck, P. (2004a). Two-mode clustering methods: A structured overview. *Statistical Methods in Medical Research*, 12, 363-394.
- Van Mechelen, I., Bock, H.-H. and De Boeck, P. (2004b). Two-mode clustering methods. In: Everitt, David Howell (Eds.), *Encyclopedia Statistics in Behavioral Science*. N.Y., Wiley.