

Classification and Data Analysis 2007

Book of Short Papers

Meeting of the Classification and Data Analysis Group of the Italian Statistical Society

eum > economia > statistica

eum > economia > statistica

The Sixth Scientific Meeting of the CLAssification and Data Analysis Group (CLADAG) of the Italian Statistical Society was held in Macerata, September 12th-14th, 2007.

Conference web site: http://cladag2007.unimc.it

The present book contains the short papers presented during this meeting.

eum edizioni università di macerata



ISBN 978-88-6056-020-9



€ 40,00

Book of Short Papers Meeting of the CLAssification and Data Analysis Group of the Italian Statistical Society



MACERATA - SEPTEMBER, 12TH-14TH, 2007

eum

Isbn 978-88-6056-020-9
Prima edizione: settembre 2007
© 2007 eum edizioni università di macerata
Vicolo Tornabuoni, 58 - 62100 Macerata
info.ceum@unimc.it
http://ceum.unimc.it
Realizzazione e distribuzione:
Quodlibet società cooperativa
Via S. Maria della Porta, 43 - 62100 Macerata
www.quodlibet.it
Stampa: Grafica Edirice Romana s.r.l., Roma

This book was realized with the financial support of

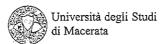






Table of contents

Preface		19
K	EY-NOTE PAPERS	
S. R. Masera, G. Mazzoni Actuarial and Continuous Time Ris financial paradigms in the new ent	sk Models: Towards a Synthesis, Changing erprise economy	23
G. McLachlan Clustering of High-Dimensional a	and Correlated Data	29
A. Montanari Classification by mixture and later	nt variable models	37
A. Rizzi Statistical Methods for Cryptograp	phy	43
R. Zang, H. Bozdogan Model Selection in Relevance Vector Complexity and Genetic Algorithm	tor Machines (RVMs) Using Information n (GA)	51
1	INVITED PAPERS	
Knowledge ext	pecialized Session 1 raction from temporal data models Organizer: D. Piccolo	
R. Baragona, S. Vitrano Statistical and numerical algorithm	ms for time series classification	65
M. Corduas Comparing time series: shape-bas	ed or structural similarities?	69
G. Scepi, G. Milone Temporal Data Mining: clustering	r methods and algorithms	73

Specialized Session 2

Statistical models with errors-in-covariates
Organizer: A. Pastore

M. Battauz, R. Bellio	
Structural analysis of linear mixed models with measurement error	79
A. Guolo Measurement error correction to I	
Measurement error correction techniques in case-control studies	83
P. Mantovan, A. Pastore	
Dynamic regression with covariate errors by covariate local clusters	87
Specialized Session 3	
Multivariate Analysis for Microarray Data Organizer: M Vichi	
M. Alfö, F. Martella, M. Vichi Biclustering of microarray data	93
A. M. Mineo, L. Augugliaro, C. Fede, M. Ruggieri A Statistical Calibration Method Based on Non-Linear Mixed Model for Affymetrix Probe Level Data	
E. Wit, G. Green	97
Random effects modelling for multivariate data from cDNA microarrays	101
Specialized Session 4	
Cluster Analysis of complex data Organizer: R. Verde	
A. Irpino, R. Verde	
Clustering linear models using Wasserstein distance	10-
F. Palumbo, A. Iodice D'Enza	107
A two-step iterative procedure for clustering of binary sequnces	111
A. Petrucci, C. T. Brownless	
Spatial Clustering Methods for the Detection of Homogenous Areas	115
Specialized Session 5	
Educational Processes Assesment by means of Latent Variables Models Organizer: S. Mignani	
M. Battauz , E. Gori	
Educational assessment in presence of heteroscedastic measurement errors	101

B. Chiandotto, F. Polverini, B. Bertaccini The effectiveness of university education: a structural equation models	125
W. J. van Der Linden Response Times on Test Items: Models and Applications	129
Specialized Session 6 Classification of complex dața Organizer: D.A. Zighed	
H. Azzag, C. Guinot, G. Venturini On data clustering with bio-inspired algorithms	135
F. de Carvalho, Y. Lechevallier, R. Verde Clustering approach on interval data	139
B. Fichet Around the ultrametric sandwich	143
D. A. Zighed Separability of classes in a multidimensional space	147
Specialized Session 7 Multidimensional Scaling Organizer: A. Okada	
G. Bove Models for asymmetry in proximity data	153
T. Imaizumi On extracting a local structure of asymmetric matrix	157
A. Okada Two-Dimensional Centrality of Asymetric Social Network	161
Specialized Session 8 Statistical Models for Public Policies Organizer: S. Ingrassia	
M. Caserta Enlarging internal regional markets: measuring the effects on local development	167
R. Innocenti, A. Giommi, C. Brownless A New Statistical Zoning for the Municipality of Firenze	. 171
M. Riani, R. Lagomarsini, A. Micozzi Robust Clustering for Performance Evaluation	175

S. Vittadini, S.C. Minotti, M. Sanarico Cluster-Weighted Modeling for Evaluating the Relative Effectiveness of Healthcare Institutions	179
Specialized Session 9 Classification models for enterprise risk management Organizer: P. Giudici	
E. Bonafede, P. Cerchiello A proposal to fuzzify categorical variables in operational risk management	185
C. Comalba Statistical models to measure IT operational risks	189
D. Fantazzini, S. Figini Predictive dynamic models for SMEs	193
Specialized Session 10 Model based clustering Organizer: D. Vicari	
C. Biernacki, A. Lourme Simultaneous Model-Based Clustering of Data Arising from Different Populations	199
G. Galimberti, G. Soffritti Multiple cluster structures and mixture models: recent developments for multilevel data	203
D. Vicari, M. Vichi Multivariate regression model with clustering of objects and variables	207
CONTRIBUTED PAPERS	
Solicited Papers Session 1 Applications in Classification and Segmentation Organizer: J. Antoch	
J. Antoch Testing the Difference of the ROC Curves in Bigamma Model	217
C. Conversano, E. Dusseldorp Classification Trunk Approach for Variable Selection and Threshold Interactions	219

	J. Klaschka Tree-based classification of EEG spectra	223
	R. Miele An Ant Colony-based Algorithm for Classification and Regression Trees"	227
	R. Siciliano, V. A. Tutore, M. Aria 3Way Trees	231
	Solicited Papers Session 2 Classification Issues in Social Network Analysis Organizer: G. Giordano, M. P. Vitale	
	V. Batagelj Clustering in Networks	237
	G. Giordano, M. P. Vitale Local Factorial Analysis and Contiguity Constraints in Social Networks	241
	Y. H. Said, E. J. Wegman, W. Sharabati, J. T. Rigsby Social Networks of Author-Coauthor Relationships	245
	S. Zaccarin, G. Rivellini Modelling network data: an Introduction to Exponential Random Graph Models	249
	Solicited Papers Session 3 Metainformation and knowledge extraction from textual data bases Organizer: S. Balbi	
	S. Bolasco, P.Pavone Automatic dictionary and rule-based systems for extracting information from text	255
	A. Canzonetti Semantic classification and cooccurrences: a method for the rules production for the information extraction from textual data	259
	F. della Ratta, B. Lorè, G. La Rocca Textual analysis perspectives on categorisation of activities in Istat survey on occupations	263
	G. Infante, M. Misuraca Text Mining Strategies for Analyzing Semi-structured Corpora	267
•		

•

•

Solicited Papers Session 4
Customer Relationship Management
Organizer: C. Davino

F. Camillo, C. Liberati A micro-data mining approach for qualitative-emotional marketing using neuro-information	273
C. Davino, R. Del Gobbo Structural Neural Networks for Modeling Customer Satisfaction	277
S. Figini, A. Roccato How to improve predictive models for database marketing applications	281
Solicited Papers Session 5 Missing Data Organizer: A. Plaia	
M. Aria, A. D'Ambrosio, R. Siciliano Robust Incremental Trees for Missing Data Imputation and Data Fusion	287
A. Plaia, A. L. Bondì Regression imputation for Space-Time datasets with missing values	291
I. Sulis, M. Porcu A multiple imputation approach in a survey on university teaching evaluation	295
Solicited Papers Session 6 Statistical methods in decision-making Organizer: R. Amenta	
G. Adelfio, M. Chiodi, D. Luzio An algorythm for earthquakes clustering based on maximum likelihood	301
F. Di Salvo, N. Ferotti, S. Consagra Hospital performance comparison: assessing inappropriate stay in the hospitals of Palermo	305
C. Gaetan, L. Greco Weighted Likelihood Inference in Gaussian Spatial Linear Models	309
Solicited Papers Session 7 Multivariate methods in social research Organizer: D. F. Iezzi	
M. Civardi, F. Crippa From University to Work: a measure of coherence between education and job	315

P. Costantini Analyzing learning effects through Latent Growth Models	319
L. Fabbris Dimensionality of scores obtained with a pair-comparison tournament system of questionnaire items	323
D. F. Iezzi, M. Grisoli Using Rasch measurement to assess the role of the traditional family in Italy	327
Solicited Papers Session 8 Statistical Models with Latent Features Organizer: L. Tardella	
S. Arima, L. Tardella A Bayesian Approach to Peabody Picture Vocabulary Test Revised	333
F. Bartolucci, A. Farcomeni Dynamic logit models for panel data based on a latent Markov heterogeneity structure	337
G. Petris, S. Petrone Bayesian inference for dynamic linear models with random variances	341
R. Rocci, A. Maruotti An INDSCAL based mixture model	345
Solicited Papers Session 9 Classificaion Trees Organizer: R. Siciliano, M. Aria	
C. Conversano, F. Mola Sequential Automatic Search of a Subset of Classifiers for Multi-Attribute Response Prediction	351
G. Galimberti, M. Pillati, G. Soffritti Comparing strategies for robust regression tree construction	355
V. A.Tutore, A. Mooijaart Optimal Scaling Trees	359
M. Vezzoli, C. Stone CRAGGING	363

•

Solicited Papers Session 10

Exploring data structures with the Forward Search
Organizer: A. Cerioli

Robust fuzzy classification		369
D. G. Calò Outlier detection via Forward Search: a propos	al based on mixture models	373
D. Perrotta, F. Torti Detecting price outliers in European trade data methods	with backward and forward	377
N. Solaro, M. Pagani The Forward Search for Metric Multidimension	al Scaling	381
Contributed Pape Clustering in Marketing Rese		
R. Menichelli, F. Palumbo Clustering of binary data in AR mining: a comp strategies	parison among different	387
I. Morlini, S. Zani Comparing approaches for clustering mixed mo marketing research	ode data: an application in	391
G. Schoier, B. Bato A modification of the DBSCAN Algorithm in a SApproach	Spatial Data Mining	395
A. Tarsitano, I. Bonafine Distance function for mixed type data		399
A. F. X. Wilhelm Association rule mining of multimedia content		403
Contributed Pape	rs Session 2	
Data Analysis in	ı Finance	
A. Amendola, M. Niglio, C. Vitale Temporal aggregation and closure of VARMA n	nodels. Some new results	409
L. Attardi, D. Vistocco Evaluating financial portfolio activeness throug	zh style analysis models	413

	R. Castellano, L. Scaccia Bayesian hidden Markov models for financial data	417
	S. Cieply, R. Abdesselam Interrelations betweeen Asset and Financial Structures: the case of French NTIC firms	421
	F. Domma, S. Giordano, P. F. Perri Modelling the dependence structure of financial assets using copula functions	425
	M. Marozzi, L. Santamaria How to Make a Complex Corporate Finance Indicator Simpler	429
· .	Contributed Papers Session 3 Mixture and Latent Class model	
	M. Alfò, A. Maruotti A Hierarchical Model for Time Dependent Longitudinal Data	435
	P. Coretto, C. Hennig Choice of the improper density in robust improper ML for finite normal mixture	439
	G. E. Montanari, M. G. Ranalli, P. Eusebi Multilevel Latent Class Models for evaluation of long term care facilities	443
	R. Piccarreta, M. Bonetti, S. Venturini Clustering-based measurement of dependence	447
	L. Scrucca Visualization of model-based clustering structures	451
	K. Yamaguchi, M. Watanabe Mixture models for multivariate count data	455
	Contributed Papers Session 4 Economics	
	V. Adomo, C. Bernini, G. Pellegrini Comparing matching methods in policy evaluation	461
	C. Cappelli, F. Di Iorio Structural breaks versus long memory, a simulation study	465
	D. De Vos, K. Vanhoof, H. Van Landeghem Symbolic Data As Measures For Logistical Performances	469
	C. Giusti Multiple imputation of income for the Labour Force Survey	473

Contributed Papers Session 5

Latent Variables

G. Adelfio Earthquakes clustering interpretation based on a second-order diagnostic approach	479
S. Bianconcini, S. Cagnone A Multilevel Latent Variable Model for Multidimensional Longitudinal Data	483
E. Otranto A time varying Hidden Markov Model with latent information	487
Contributed Papers Session 6 Non Parametric Regression	
D. Coin A test for Normality with high Power against Symmetrical Alternatives	493
L. De Benedictis, M. Callegati, M. Tamberi Semiparametric analysis of the specialization-income relationship	497
R. Miglio Non Parametric Procedures to Explore Interactions in Epidemiological Studies	501
A. Zirilli, A. Alibrandi SMART Regression: an application to multidimensional data in a medical context	505
Contributed Papers Session 7 Cluster Structures	
A. Alibrandi, M. Giacalone A Comparison among Statistical Criteria to Establish the Number of Clusters in Multidimensional Data	511
A. Balzanella, E. Romano, R. Verde Outliers detection for a Robust Curves Clustering Algorithm	515
A. Iodice D'Enza Visual Monitoring Tool of Association Patterns in Binary Data Flows	519
G. Notarstefano, R. Scuderi Polarization and Dynamic Clustering: A Review	523
S. Polettini, L. Di Consiglio Disclosure risk estimation with survey data: a comparative study	527

Contributed Papers Session 8

Customer Satisfaction

L. Berzieri Demographic snapshot of Parma electoral colleges	533
P. Chirico Common Optimal Scaling for Customer Satisfaction Multidimensional Models	537
M. Civardi, E. Zavarrone Ordinal Scales and Customer Satisfaction: an Applet	541
P. Ferrari, C. V. Fiorio, L. Pagani A two-step procedure to analyze users' satisfaction	545
Contributed Papers Session 9 Discrimination	
R. Abdesselam Mixed Discriminant Analysis	551
E. Brichieri Colombi, C. T. Brownlees, G. M. Gallo Predicting Binary Outcomes Using Flexible and Parsee's Methods	555
S. Ingrassia, I. Morlini Computational experiences with the equivalent number of degrees of freedom of neural networks	559
R. Sanzo, E. Trinca A multivariate approach to non-response treatment for Foreign control companies Survey	563
Contributed Papers Session 10 Robust Classification	
A. Corbellini, P. Ganugi, L. Crosato Robust sequential Fitting of the Pareto II distribution: Theoretical and Computational Aspects	569
M. Daniele, A. Bondì Outliers in hierarchical models: an application to air pollution data	573
A. Gottard, S. Pacillio Robustifying SINful procedure	577
P. Mascia, R. Miele, F. Mola Robustness Issues in Classification and Regression Trees	581
S. Scippacercola The Progressive Single Linkage Algorithm Based on Minkowski Ultrametrics	585

Contributed Papers Session 11

Bayesian Network

A. Brogini, D. Slanzi Several computational studies about variable selection for Bayesian networks	591
M. Mascherini Exploring eBusiness Readiness Indicators using Bayesian Networks	59:
L. Scaccia, F. Bartolucci , Bayesian inference for marginal models under equality and inequality constraints	599
Contributed Papers Session 12 Textual data	
M. Dimai, N. Torelli Clustering textual data by Latent Dirichlet Allocation: Applications and Extensions to Hierarchical Data	605
D. F. Iezzi, M. Giorgino Intimate femicide in Italy: a model to classify how killings happened	609
S. Volo, C. Fisichella Evaluating tourists experiences	613
Contributed Papers Session 13 Data Analysis	
M. Greenacre Power Transformations in Correspondence Analysis and Related Methods	619
G. Menardi, N. Torelli Multimodal projection pursuit using the adjusted critical bandwidth	623
R. Romano, T. Naes, P. Brockhoff The use of analysis of variance and three-way factor analysis methods for studying the quality of a sensory panel	627
POSTER SESSION	
R. Belicchi, F. Crippa, L. Cesana, R. Daini Lezak TinkerToy Test: freely making constructions to assess executive functions. An Italian normative sample	633

	A. Belliggiano, S. Staffieri Private label organic food products: profile and behaviour of Italian consumer characterized through multivariate approach
• •	G. Boari, G. Cantaluppi, A. De Lauri Handling missing values in PLS path modelling: comparison of alternative procedures
	S. Bozza, F. Taroni, F.Guy, M. Schmittbuhl Taxonomic membership in great apes: a probabilistic approach from mandibular morphology
	S. De Cantis, M. Ferrante Data quality and errors correction in Italian official statistics on guests of accommodation establishments
	C. Fukae The Myths of Global Warming and Nuclear Power: Evidence from public opinion survey in Japan
	R. Gismondi, A. R. Giorgi, T. Pichiorri Reducing Sampling Error using Deterministic Sample Selection in the Italian Retail Trade Monthly Survey
	R. Gismondi, M. A. Russo Selection and Statistical Treatment of Variables for Evaluating Quality of Life in the Italian Provinces
	U. Magagnoli, P. M. Chiodini Unilateral Hypothesis Tests of Efficiency Functions with Heteroscedastic Errors: an Iterative Procedure
	A. M. Olivieri Excursionist flows and the attractiveness on the Aeolian Islands
	A. G. Quaranta, S. Fedeli, E. Voltattorni Basel 2: firm's Scoring via Cluster Analysis and Artificial Neural Networks
	Ph.D. SPECIAL SESSION
	G. Adelfio Point processes residual analysis and asymptotical distribution of transformed versions of some second-order statistics
	V. Adorno Program evaluation with continuous treatment: theoretical considerations and empirical application

Health related Quality of Life: some methodological aspects of Rasch analysis	689
P. Eusebi Many Facets Rasch Models: a Hierarchical Latent Class approach	693
M. Matteucci A multidimensional Item Response Theory approach for the University guidance	697
R. Ragona A State-Space Model for the Generation of Spatial Maps	701
E. Romano A new strategy for curves clustering based on free knots spline estimation	705
I. Sulis Measuring students' assessment on 'university course quality' using mixed-effects models	709
Indice degli autori	713

it genera to ies of false hat is cases unsatisfac-

Data quality and errors correction in Italian official statistics on guests of accommodation establishments (1)

Stefano De Cantis, Mauro Ferrante
Dipartimento di Metodi Quantitativi per le Scienze Umane
Università degli Studi di Palermo
Viale delle Scienze – Edificio 13, 90133 Palermo
decantis@unipa.it, mauroferrante@unipa.it

Abstract: This paper discusses the quality of data in Italian survey on guests of accommodation establishments by looking at some basic conditions of reliability of the elementary records. Monthly data referred to the 32 Sicilian tourism promotion agencies for seven years are checked and the main kinds of error are presented. Some hypothesis with reference to the sources of the errors and to their relative correction procedures are discussed. The method proposed, relatively simple to apply, could be implemented as a part of the standard data pre-processing.

Keywords: Tourism statistics, error rate, reliability

1. Introduction

One of the most important official Italian survey on tourism is the so-called Survey on guests of accommodation establishments (ISTAT, 2004). It is usual to refer to these data as information on tourists' arrivals (ISTAT, 2006, p.209). This is a conceptual mistake for several aspects. First of all because information about characteristics of guests are not available, this makes impossible to distinguish tourists from other guests (European Commission, 2007, Technical notes). Secondly, there is a huge number of accommodation establishments (such as private accommodations), not registered to the Italian Commercial Facilities Register for which information on guests are not covered by this survey. Moreover, the official accommodation establishments cannot declare the real number of arrivals and/or of presences, to avoid direct or indirect taxation. Finally, and even more important, a single guest determines more than one arrival on a given destination for a given month, if he/she changes accommodation establishment during his/her trip, generating errors resulting in an overestimate of the number of tourists' arrivals, and consequently in an underestimate of the average length of stay (Vaccina and Parroco, 2006). This problem is well documented but, unfortunately, there is no direct way to improve the reliability of the survey on measuring the number of tourists without changing the framework of the survey (for some pioneers researches on a single local district, see Parroco and Vaccina, 2006). In this paper we investigate the way to improve the quality of data controlling for several basic conditions, by checking the internal coherence of records.

⁽¹⁾ The present paper is a common work and responsibility of both authors, however S. De Cantis wrote sections 1 and 3, and M. Ferrante wrote sections 2 and 4.

2. Detecting errors

The arrivals (number of guests spending at least one night) and the presences (number of nights spent by guests) in accommodation establishments are recorded by Italian tourism promotion agencies, according to the so-called Istat CTT/1 model (ISTAT, 2004). It is constituted by four main sections derived from the combination of two factors: region of origin of guests, and typology of accommodation used. According to our framework of analysis, the elementary record is composed by a pair of values as follows:

$$(x_{t,w,i,j}, y_{t,w,i,j})$$
 with $x_{t,w,i,j}$ and $y_{t,w,i,j} \in \mathbb{N}^+ \ \forall t, w, i, j$ (1)

where: x and y are the arrivals and the presences on accommodation establishments; t (time) ranges from 1 (jan-1999) to 84 (dec-2005); w represents the 32 Sicilian tourism promotion agencies; t indicates the country (if foreign) or the region (if Italian) of origin of guests, and it ranges from 1 to 76 (21 Italian regions and 55 foreign countries); and t represents the ten typologies of accommodation establishments, resulting in a total number of pairs equals to 2,042,880. Since one arrival generates at least one presence, it holds:

$$x_{t,w,i,j} \le y_{t,w,i,j} \quad \forall t, w, i, j \tag{2}$$

If we define the average length of stay (z) as the ratio between presences and arrivals, the condition (2) becomes:

$$z_{t,w,i,j} = y_{t,w,i,j} / x_{t,w,i,j} \ge 1 \qquad \forall t, w, i, j$$
(3)

To identify all the possible violations of the condition (3), monthly pairs of arrivals and presences referred to Sicily from 1999 to 2005 were classified according to eight possible conditions, as shown in Table 1. If we exclude records referred to condition A (null pairs), only the records corresponding to the conditions B and C can be considered as correct. On the other side, records falling into conditions from D to H must be considered as incorrect records. As result, we found a raw error rate equals to 2.4%. It should be noted that we also found several records falling into condition B, that are "pathological". For example, 1,117 (0.3%) records have an average length of stay greater than 60 days.

Table 1: Number of correct and incorrect records according to different conditions on values of arrivals and presences in Sicily (1999-2005).

	Conditions	on		Number of pairs	%	%
	Arrivals	Presences	Average length of stay		,	
Α	$x_{t,w,i,j}=0$	$y_{t,w,i,j} = 0$	Indefined	1,602,035		78.4%
В	$x_{t,w,i,j} > 0$	$y_{t,w,t,j} > 0$	$ y_{t,w,i,j}/x_{t,w,i,j}>1$	347,985	78.9%	
	$x_{t,w_i,i,j} > 0$	$y_{t,w,i,i} > 0$	$y_{t,w,i,j}/x_{t,w,i,j}=1$	82,293	18.8%	
D	$x_{t,w_i,i,j}=0$	$y_{t,w,i,j} > 0$	$y_{t,w,i,j}/x_{t,w,i,j} \to +\infty$	9,764	2.2%	
E	$x_{t,w,i,j} > 0$	$y_{t,w,i,j}=0$	$y_{t,w,i,j}/x_{t,w,i,j}=0$	252	0.1%	
F	$x_{t,w,i,j} > 0$	$y_{t,w,i,j} > 0$	$ y_{t,w,i,j}/x_{t,w,i,j} \leq 1$	548	0.1%	
	$x_{i,w,i,j} > 0$	$y_{t,w,i,j} < 0$	$ y_{t,w,i,j}/x_{t,w,i,j} > 1$	3	0.0%	
H	Other cond			0	0.0%	
		Subto	440,845	100.0%	21.6%	
	-		Total	2,042,880		100.0%

Source: Osservatorio turistico Regione Siciliana, and authors' calculations.

To include all the possible cases, it is possible to explicit the residual conditions H, as follows:

- H1: $x_{t,w,i,j} > 0$ and $y_{t,w,i,j} < 0$ and $|y_{t,w,i,j}/x_{t,w,i,j}| < 1$;
- H2: $x_{t,w,i,j} > 0$ and $y_{t,w,i,j} < 0$ and $|y_{t,w,i,j}/x_{t,w,i,j}| = 1$;
- H3: $x_{t,w,i,j} < 0$ and $y_{t,w,i,j} < 0$ and $|y_{t,w,i,j}| < 1$;
- H4: $x_{t,w,i,j} < 0$ and $y_{t,w,i,j} < 0$ and $|y_{t,w,i,j}| = 1$;
- H5: $x_{t,w,i,j} < 0$ and $y_{t,w,i,j} < 0$ and $|y_{t,w,i,j}/x_{t,w,i,j}| < 1$;
- H6: $x_{t,w,i,j} = 0$ and $y_{t,w,i,j} < 0$ and $|y_{t,w,i,j}/x_{t,w,i,j}| \rightarrow +\infty$
- H7: $x_{t,w,i,j} < 0$ and $y_{t,w,i,j} > 0$ and $|y_{t,w,i,j}/x_{t,w,i,j}| < 1$;
- H8: $x_{t,w,i,j} < 0$ and $y_{t,w,i,j} > 0$ and $|y_{t,w,i,j}| = 1$;
- H9: $x_{t,w,i,j} < 0$ and $y_{t,w,i,j} > 0$ and $|y_{t,w,i,j}| > 1$.

nber of ourism). It is gion of vork of

(1) nents; t ourism rigin of ; and j

3. Describing the sources of error variability

The above classification is useful to determine the possible sources of error, and subsequently the more appropriate ways to correct the data, as it will be discussed into the last section. Since the distribution of errors on records resulted to be not uniform according to several dimensions of our data, in Table 2 we reported the raw error rates specific by tourism district agency, by months, by years, by typology of accommodation establishment, respectively. Moreover, to consider the joint effect of these variables we estimated, only for descriptive purposes, several logit models. Let $n_{l,m,v,w,i,j}$ and $n_{2,m,v,w,i,j}$ the number of incorrect and correct records respectively, for the m-th month, for the v-th year, for the w-th tourism district, for the i-th nationality (1 = Italian, 2 = Foreign), and for the j-th accommodation establishment. The logit model for the probability of incorrect records — with a simple structure, and the best goodness of fit (interactions between factors did not appeared interesting) — can be expressed as follows:

$$\log(n_{1,m,v,w,i,j}/n_{2,m,v,w,i,j}) = \log(n_{1,m,v,w,i,j}) = \theta + \theta_m^M + \theta_v^V + \theta_w^W + \theta_i^J + \theta_j^J$$
(4)

 $\forall m=1,2,...,12;\ v=1,2,...,7;\ w=1,2,...,32;\ i=1,2;\ j=1,2,...,10$ and $\theta_{12}^{w}=\theta_{1}^{v}=\theta_{12}^{v}=\theta_{2}^{i}=\theta_{3}^{i}=0$. Through the model (4) the effect of each factor is evaluated controlling for the other ones, and Table 2 reports the corresponding parameters (Odds Ratios, OR). However, the results are almost comparable. We found a different raw error rate between records related to Italian guests and foreign ones (1.9% vs 2.8%, respectively), and a risk of incorrect record almost double in the case of Italian guests, compared to the foreign ones. A large variation appears in errors rate and in odds ratios, with reference to different tourism districts. Agrigento resulted to be the district with the lowest associated risk (adj. OR = 0.02, rank 32), considering as reference (OR = 1) the district of Trapani (others municipalities).

Table 2: Raw error rates and adjusted odds ratios for incorrect records

Table 2. Raw error rates and adjusted odds ratios for incorrect records.																																
Destination		Agrigento others	- 1	Gela	ſ	Caltanissetta					Calanta (others)		Piazza		Capo D'Orlando	Glardini	Isole Eolie	Messina	Afilazzo	Patti	Тооттіпа	Alessina	Cefalir	Patermo	Palermo	Rogura	Ragusa	Siracusa	Stracusa	Erice	Trapani	Frapani (others)
Error rate Adj. OR	0.02	1.76	3.21	4.25	0,59	2,32	0.07	0.88	1.50	1.20	4,58	1.05	0.94	5.07	2.79	2.46	3,37	6.77	1,55	2.11	2.13	2.37	2 80	2 69	2 41		1	1			- 1	
Rank	37	21	5.14	4.44	0.57	2.07	0.05	0.64	0.67	0,48	2.94				2.25	2.13	1.89	6.77 4.17	1.24	1.61	1,17	1.65	1.58	1.24	1.45	0.04	0.10	1 90	3.13	1 70	1.55	1.89
	, V-4		4		20		30	25	.24	28	5	23	26	2	6	7	9	3	18	14	20	13	15	18	16	30	29	9	11	17	17	1.00
Month		- 1	Jar	7	Fe	h	$\overline{}$	Mar		A		—r	Ма		_	-			_			_			_			_	<u> </u>		111	<u></u>
Error r	rate	T	4.7	$\overline{}$.59	~├	2.1				+				un		Jul	_	1	lug		Se	р.		_Oc	f	7	Ιον	П	Dec	_
Adj. O		-	1.9				- -			—~	.67			73	_	1.89		1.	82		2.2	6	2	.25	T	2.	90	1	3.48		2.	77
	к.	-	1.9	3	_	.99	-4-	0.7	6	(1.53	_1	_0.	54		0.51	7	0.	52	Т	0.6	6	- 0	.67	7	0.		→—	1.31			00
Rank		_L	1			4	_L	6		!	11		1	0	T	9	\neg	1	2	+	8	_	_	7	-		/ -	+	_			
**																			_	I				_	_	<u> </u>	,		2		3	i _
Year			_ _	_	999		1	2	200		Т		200	1	\neg		20	02	_		2	003		$\overline{}$		20	1.7		_			
Error r	ate			2	2.87		-	2	2.62		2,46					2.43				┿				2004					$\overline{}$	2005		
Adj. OF	R			1	.46		\Box	$\overline{}$.35	_	+	1.23			-				 	2.01			4	2.15				2.		2.40		
Rank		_		_	1	_	_		2	_	+	_		_	\dashv	1.17			0.90) 0		0.92					1.00			
					•		느.		<u>-</u>			_	3	_	[4			1		7		- [-	6			5		_		

Accommodation			*	Hotels				Camping	Rural	
establishment	5	4	3	2		<u> </u>	Rented	&	tourism.	Other
typology Error rate	stars	stars	stars	stars	1 star	Residences	facilities	resorts	facilities	establishments
Adi. OR	1.66	1.28	1.37	1.67	1.96	10.15	3.59	3.15	2.25	4.91
Rank	1.03	0.98	1.00	1.20	1.40	11.53	3.48	3.18	2.08	4.63
		10	_9	<u> </u>	7	1	3	4	5	<u> </u>

Error rates demonstrate also a dependence on months, with relatively lower values in high season months (e.g. Apr.-Sep.), and with higher values in low season months (e.g. Nov.-

Feb.). Instead, during the years considered, it seems not to be strong differences, both in the error rates and in the odds ratios. Finally, considering different accommodation establishments, it is possible to highlight a strong difference between the relatively small raw rates of hotels from 5 to 1 star, and the other typologies; with a maximum in correspondence of *Residences* (Error rate = 10.15%, OR = 11.53).

4. Considerations on errors correction and conclusions

Recalling the conditions expressed in Table 1, it is possible to discuss some considerations related to the causes of errors and to their relative correction procedure. Conditions D and E are referred to the situation in which one of the two elementary data (arrivals or presences) is zero. In these cases, considering the zero value as a missing one, a procedure of imputation can be based on an "appropriate value" of the average length of stay as an indicator of the relationship between the existing value and the missing one. Condition F is not easy to solve in term of imputation procedure. We assume that these errors derive from a mixture of correct and incorrect (deriving, for example, from condition D) records. In this case if we suppose that the most reliable data is the one related to the arrivals, we will use, also in this case, an "appropriate value" of the average length of stay, to derive the correspondent value of presences. Finally, condition G can reveal a simple mistake on editing; thus a procedure of correction can consists in a simple substitution of the value of presences with its absolute value. With reference to the remaining conditions included in H, we could suggest the same procedure used for G, with reference to the sub-cases from H2 to H4, and H8-H9, while with reference to the other possible situation, here not discussed in detail, a combination of the procedures used for D, E, and F, could be adapted. In conclusion, this paper showed how to control for reliability of the elementary records of Italian survey on guests of accommodation establishments, by checking jointly data on arrivals and presences. The results showed an absence of accurate controls on the data sent by the accommodation establishment' responsible to the local tourism promotion agency. The method proposed, relatively simple to apply, could be implemented as a part of the standard data pre-processing. This control could be improved through the use of other constraints on the supply side derived, for example, from the total number of bed places available. It is important to perform these controls before aggregating the data, since the process of aggregation could determine a subsequent impossibility of detecting Moreover, since the presence of errors seems to have some systematic components, a deeper analysis could suggests the more appropriate way to eliminate the causes of errors or to impute appropriate values, improving the quality of data.

References

European Commission (2007) Tourism statistics, Luxembourg, http://ec.europa.eu/eurostat

ISTAT (2004) Statistiche del turismo: Anno 2002. Roma, www.istat.it

ISTAT (2006) Italian Statistical Abstract, 2004. Roma, www.istat.it

Parroco A.M., Vaccina F. (2006) Methodological issues and first results of a study to identify and evaluate the error of some tourism statistics. Proceedings of XLIII Scientific Reunion of Italian Statistical Society (SIS), Torino. Cleup, Padova.

Vaccina F., Parroco A.M. (2006) Referring to Space and Time when using territorial data: The case of Touristic Arrivals. Proceedings of International Statistical Institute (ISI) 55th Session, Sidney, Australia.