FROM DATA MINING TO KNOWLEDGE MINING: SYMBOLIC DATA ANALYSIS AND THE SODAS SOFTWARE

> E. Diday University of Paris IX Dauphine and INRIA

AIM

FROM HUDGE DATA IN AN ECONOMIC WAY

-Extract new knowledge

-Summarize

-Concatenate

-Solve confidentiality

-Explain correlation

HOW? By working on HIGHER LEVEL UNITS called CONCEPTS necessary described by more complex data extending Data Mining to Knowledge Mining.

OUTLINE

1) THE MAIN IDEA: FIRST AND SECOND ORDER OBJECTS.

2)THE INPUT OF A SYMBOLIC DATA ANALYSIS:

SYMBOLIC DATA TABLE.

3) SOURCE OF SYMBOLIC DATA: FROM DATA

BASES, FROM CATEGORICAL VARIABLES.

4) MAIN OUTPUT OF SYMBOLIC DATA ANALYSIS

ALGORITHMS: SYMBOLIC DESCRIPTIONS AND

SYMBOLIC OBJECTS.

5) THE MAIN STEPS OF A SDA.

6)TOOLS OF SYMBOLIC DATA ANALYSIS

7) SYNTHETICAL VIEW OF THE SODAS PROJECT

THE MAIN IDEA: FIRST AND SECOND ORDER OBJECTS THE ARISTOTLE ORGANON (IV B.C.) CLEARLY DISTINGUISHES "FIRST ORDER OBJECTS" (AS THIS HORSE OR THIS MAN) CONSIDERED AS A UNIT DESCRIBING AN INDIVIDUAL OF THE WORLD , FROM "SECOND ORDER OBJECTS" (AS A HORSE OR A MAN) ALSO TAKEN AS A UNIT DESCRIBING A CLASS OF INDIVIDUALS.



FROM FIRST ORDER OBJECTS TO SECOND ORDER OBJECTS IN OFFICIAL STATISTICS

Units	Classes	Descr. Var. of the Units			
Case nº	Region	Bedroom	Dining- Living	Socio-Econ Group	
11401	Northern- Metropolitan	2	1	1	
11402	Northern- Metropolitan	2	1	3	
11403	Northern- Metropolitan	1	3	3	
12315	East-Anglia	1	3	3	
12316	East-Anglia	2	2	1	
14524	Greater-London	1	2	3	

	IN	DFFICIAL	ST.	ATISTICS			
Classes	Descriptive	variable of	the	units			
Region	Bedroom	Dining-Liv	So	cio-Ec gr			
orthern- Aetropolitan	2	1	1				
orthern- Aetropolitan	2	1	3				
Northern- Aetropolitan	1	3	3				
ast-anglia	1	3	3				
ast-anglia	2	2	1				
last-anglia	1	2	3				
Classos	Descriptive	veriebles of	the	alassos			
Region	Bedroom	Dining-Li	v	Socio-Fe	ar		
Northern- Metropolitan	(2\3) 2, (1\3) 1	(2\3) 1, (1\3	3) 3	(1\3) 1, (2\	3) 3		
East-anglia	(2\3) 1, (1\3) 2	(2\3) 2, (1\3	3) 3	(2\3) 3, (1\3	3) 1		



EXAMPLE OF SYMBOLIC DATA TABLE

PRODUCT	WEIGHT	TOWN	COLOUR
PRODUCT 1	3.5	Londres	{red, white, yellow}
PRODUCT 2	[3,8]	{Paris, Londres }	
PRODUCT 3	{3.1 , 4.6, 7.2}		{ 0.3 red, 0.7 green}
PRODUCT 4	[(0.4) [2,3[, (0.6) [3, 8]]		

THE CELLS CAN CONTAIN:

-SEVERAL QUALITATIVE OR QUANTITATIVE WEIGHTED VALUES

-INTERVALS

- HISTOGRAMS



SOURCE OF SYMBOLIC DATA

.FROM CATEGORICAL VARIABLES: - GIVEN. (AS « TYPE OF EMPLOYMENT ») - OBTAINED BY CLUSTERING.

.FROM DATA BASES: QUERY CREATING A NEW CATEGORICAL VARIABLE: cartesian prod

.FROM EXPERT: NATIVE SYMBOLIC DATA:

Scenario of road accidents, species of insects

.FROM CONFIDENTIAL DATA IN ORDER TO HIDE THE INITIAL DATA BY LESS ACCURACY

.FROM STOCHASTIC DATA TABLE:

THE PROBABILITY DISTRIBUTION , THE HISTOGRAM THE PERCENTILES OR THE RANGE OF ANY RANDOM VARIABLE ASSOCIATED TO EACH CELL OF A DATA TABLE

EXAMPLE

	Mathematics	Physics	Litterature	
Tom	X _M	Хр	XL	
Paul				

 $\mathbf{X}_{\mathbf{M}}$ is the random variable which associates to each

exam of TOM his mark in mathematics.

From X_M several kinds of symbolic objects can be

defined by using in each cell: - Probability distr.

- Histograms

- Inter-quartile intervals



Example:Districts description					
Districts	male-full- time- employee %	male-part- time- employee %	male-self- employed %	Z	
D1	8%	5%	20%	Z1%	
D2	12%	9%	15%	Z2%	
Dn					

Example of Rule defining class

- **R1**: male-full-time-employee%(X,low) ∧ malepart-time-employee%(X,low) ∧ neighbor(X,Y) ∧ comm-activities(Y,high)
- \rightarrow male-self-employed%(X,high)
- 70 districts X satisfy the rule: the low percentage of full-time and part-time male employees in district X adjacent neighbor of *Y*, with many commercial activities, implies a high percentage of self-employed males in *X*.

Descript of rules extent	male-full- time- employee%	male-part- time- employee %	male-self- employed%	Z
R1	[8%,12%]	[5%, 9%]	[20%, 15%]	[Z1%,Z2]
R2	[2%, 6%]	[4%,8%]	[18%,14%]	

MA DAT	IN OUTPUI A ANALYSI	OF SYMBOLIC S ALGORITHMS:			
	SYMBOLIC D	DESCRIPTIONS			
	SYMBOLIC OBJECTS.				
SYN	BOLIC DE	SCRIPTIONS			
Description	AGE	SDC			
Description	AGE	SPC			
Description D1	AGE {12,20,28}	SPC {employee,worker}			







THE MEMBERSHIP FUNCTION« a » MODAL CASE

$$\begin{split} & S = (a, R, d): \\ & a(w) = [age(w) R_1 \{(0.2)12, (0.8) [20, 28]\}] \land \\ & [SPC(w) R_2 \{(0.4) employee, (0.6) worker\}] \\ & a(w) \in [0,1]. \\ & First approach: simple or flexible matching \\ & R = (R_1, R_2): r R_i q = \sum_{j=1,k} r_j q_j e^{(r_j - \min(r_j, q_j))}. \\ & Second approach: \\ & Probabilistic: if dependencies, copulas, \\ & derivation of the joint distribution, \\ & transforming the joint density in [0,1]. \end{split}$$

EXTENT OF A SYMBOLIC OBJECT S: BOOLEAN CASE: EXT(s) = {W $\in \Omega / a(W) = TRUE$ }. MODAL CASE EXT_{\alpha} (S)= EXTENT_{\alpha} (a) = {W $\in \Omega / a(W) \ge \alpha$ }.







APPLY

SYMBOLIC DATA ANALYSIS TOOLS

- Correlation, Mean, Mean Square
- Histogram of a symbolic variable
- Dissimilarities between symbolic descriptions
- Clustering of symbolic descriptions
- Principal component Analysis
- Decision Tree
- Graphical visualisation of Symbolic Objects







DE CARVALHO'S DISSIMILARITY MEASURES

Five different similarity measures s_i , i = 1, ..., 5, are defined:

s i	Comparison Function	Range	Property
S ₁	$\alpha/(\alpha+\beta+\chi)$	[0,1]	metric
S ₂	$2\alpha/(2\alpha+\beta+\chi)$	[0,1]	semi metric
S 3	$\alpha/(\alpha+2\beta+2\chi)$	[0,1]	metric
S 4	$\frac{1}{2} \left[\frac{\alpha}{(\alpha + \beta) + \alpha} \right]$	[0,1]	semi metric
S₅	$\alpha/[(\alpha+\beta)(\alpha+\chi)]^{\frac{1}{2}}$	[0,1]	semi metric

The corresponding dissimilarities are $d_i = 1 - s_i$. The d_i are aggregated on p variables by the generalised Minkowski metric, thus obtaining: SO 1

$$d^{i}(a,b) = \sqrt[q]{\sum_{j=1}^{p} \left[w_{j} d_{i}(A_{j},B_{j}) \right]^{q}} \qquad 1 \le i \le 5$$





























Some recent advances:	
- Mixture decomposition of Distributions of distributions (by Copulas, Dirichlet and Kraft stochastic process)	
- Stochastic Symbolic Conceptual lattices using capacity theory	
- Symbolic class description	
-Symbolic Regression	
-NEXT FUTUR	
Spatial symbolic clustering by pyramids	
- Symbolic time series.	
- Consensus between different description of the same set of units	

AIM ATTAINED FROM HUDGE DATA BASES IN AN ECONOMIC WAY WE ARE ABLE TO: -Extract new knowledge -Summarize -Concatenate -Solve confidentiality -Explain Correlation HOW? By working on HIGHER LEVEL UNITS extending Data Mining to Knowledge Mining.

CONCLUSION

Symbolic Data Analysis is an extension of standard data analysis therefore

First principle: any Symbolic Data Mining method must have as a special case method of Data Mining on standard data.

Second principle : the output must be a symbolic description or symbolic object

New problems appear as the quality, robustness and reliability of the approximation of a concept by a symbolic object, the symbolic description of a class, the consensus between symbolic descriptions etc..

Due to the intensive development of the information technology the great chapters of the standard statistics will have to be think in these new terms.

References

SPRINGER, 2000 :

"Analysis of Symbolic Data"

H.H., Bock, E. Diday, Editors . 450 pages.

JASA (Journal of the American Statistical Association) "From the Statistic of Data to the Statistic of Knowledge: Symbolic Data Analysis" Billard, Diday June, 2003.

Electronic Journal of S. D. A.: ESDA E. Diday, R. Verde, Y. Lechevallier

Download SODAS and SODAS information : www.ceremade.dauphine.fr/~touati/sodas-pagegarde.htm

THE SODAS 2 SOFTWARE FROM ASSO



