

MEASURING SIMILARITY OF COMPLEX AND HETEROGENEOUS DATA IN CLUSTERING OF LARGE DATA SETS

Helena Bacelar-Nicolau¹, Fernando Nicolau², A.Urea Sousa³, Leonor Bacelar-Nicolau⁴

¹*Universidade de Lisboa. FPCE, Laboratório de Estatística e Análise de Dados,
Lisboa, Portugal*

²*Universidade Nova de Lisboa, FCT, Departamento de Matemática, Caparica, Portugal*

³*Universidade dos Açores, Departamento de Matemática, Ponta Delgada, Açores, Portugal*

⁴*Universidade de Lisboa. FML, Laboratório de Biomatemática & Instituto de Medicina
Preventiva. Lisboa, Portugal*

Abstract:

Cluster analysis or classification usually concerns a set of exploratory multivariate data analysis methods and techniques for finding a clustering structure on a dataset. That may refer either to groups of statistical data units or to groups of variables. In this work we deal with a generalization of this paradigm concerning clustering of complex data described by three different types of variables, frequently present in a three-way context. We obtain compatible versions of the same affinity coefficient for measuring similarity between statistical data units described by those three types of variables. A global generalized similarity coefficient is analyzed for such kind of mixed data, often arising in data mining or knowledge mining.

Keywords: cluster analysis, different type variables, similarity coefficient, three-way data