# Machaon CVE: Cluster validation for gene expression data

*Nadia Bolshakova[a,*] and Francisco Azuaje[b]*

[a]*Department of Computer Science, Trinity College Dublin, Dublin 2, Ireland*
[b]*School of Computing and Mathematics, University of Ulster, Jordanstown Co. Antrim, BT37 0QB Northern Ireland, U.K*

[*]To whom correspondence should be addressed

## ABSTRACT

**Summary:** This paper presents a cluster validation tool for gene expression data. *Machaon-CVE* (Clustering and Validation Environment) system aims to partition samples or genes into groups characterised by similar expression patterns, and to evaluate the quality of the clusters obtained.

**Availability:** The program is freely available for non-profit use on request at http://www.cs.tcd.ie/Nadia.Bolshakova/Machaon.html

**Contact:** Nadia.Bolshakova@cs.tcd.ie

**Supplementary information**: http://www.cs.tcd.ie/Nadia.Bolshakova/Machaon.html

## INTRODUCTION

DNA microarray technologies allow measuring the expression of thousands of genes in parallel under multiple experimental conditions. Genomic and post-genomic studies (Schena *et al*., 1995), such as disease diagnosis, drug discovery and toxicological research have been benefited from it (Debouck and Goodfellow, 1999, Gray *et al*., 1998).

An important step in the analysis of gene expression data is the detection of samples or genes with similar expression patterns. Several clustering algorithms have been developed for gene expression data. Also solutions to systematically evaluate the quality of the clusters have been presented (Bolshakova and Azuaje, 2002). The prediction of the correct number of clusters in a data set is a critical problem in unsupervised classification. Various cluster validity indices have been proposed to measure the quality of clustering results (Azuaje, 2002; Dudoit. and Fridlyand, 2002). Clustering algorithms may require the a-priori definition of the number of clusters. Hence, a clustering algorithm can be executed several times, with different number of clusters in each run, and the clustering partition that optimises a validity index is selected as the best partition. Previous studies have not provided integrated tools for both clustering and automatically assessing the quality of the resulting clusters. Therefore, there is the need to design software platforms, which integrate clustering and validation methods for predicting the optimal number of clusters in gene expression data analyses.

The *Machaon CVE* is a cross-platform Java-based tool, which offers multiple clustering and validity methods for DNA microarray data analysis. It aims: a) to partition samples or genes into groups characterised by similar expression patterns, and b) to evaluate the quality of the clusters obtained.

## SYSTEM OVERVIEW

The software is implemented as a multi-window Java application, which allows working with different datasets, clustering (hierarchical and k-means) and validation (C-index, Davis-Bouldin, Dunn's, Goodman-Kruskal and Silhouette indices) algorithms, and results simultaneously. For further information on the implementation, of these algorithms the reader is referred to the supplementary information page. The system supports several modifications of tabular data format widely used by third-party clustering tools (Herrero *et al.*, 2001) Moreover, an XML-based format is being designed to address some of the limitations observed in traditional formats, such as inability to store multiple clustering and validation results within a dataset.

Multiple clustering may be applied to a single dataset and the results may be easily compared. Every clustering result may be selected and validated across a number of parameterised validation methods. Both clustering and validation results are represented as two-level tree in the bottom of the corresponding data set window (Figure 1). Clustering indices are also displayed in additional columns of a data set table. Every such column is associated with a single partition. The results of a hierarchical clustering can also be displayed using dendrograms. Users may choose from a collection of clustering and validation techniques, compare the results from each method and generate interpretations.

Several methods for measuring gene-to-gene (or sample-to-sample), intercluster and intracluster distances can be used in any combination. This is important to research the influence of different distance metrics on both clustering and validation. Machaon CVE provides data normalization functionality, which may be either selected as an option of clustering/validation or used to produce a normalized dataset.

Apart from the clustering and validation results, the system shows, if known, the natural classification structure of the data (leukemia types in the example illustrated in Figure 1), which allows comparisons against clustering results and validation analyses across natural classes.

Despite the fact that *Machaon CVE* was developed for DNA microarray expression analysis applications, it may be effectively used for clustering/validating other biomedical and physical data with no limitations.

## ACKNOWLEDGEMENTS

## REFERENCES

Azuaje,F. (2002) A cluster validity framework for genome expression data, *Bioinformatics*, **18**, 319-320.

Bolshakova, N. and Azuaje,F. (2003) Cluster validation techniques for genome expression data, *Signal Processing*, **83**, 825-833.

Debouck,C. and Goodfellow,P.N. (1999) DNA microarrays in drug discovery and development, *Nature Genet.*, **21**, 48-50.

Dudoit,S. and Fridlyand,J. (2002) A prediction-based resampling method for estimation the number of cluster in a dataset, *Genome Biology*, **3**, 1-21.

Grey,N.S., Wodicka,L., Thunnissen,A.M., Norman,T.C., Kwon,S., Espinoza,F.H., Morgan,D.O., Barnes,G., LeClerc,S., Meijer,L., Kim,S.H., Lockhart,D.J. and Schultz,P.G. (1998) Exploiting chemical libraries, structure, and genomics in the search for kinase inhibitors. *Science*, **281**, 533-538.

Herrero,J., Valencia,A. and Dopazo,J. (2001) A hierarchical unsupervised growing neural network for clustering gene expression patterns, *Bioinformatics*, **17**, 126-136.

Schena,M., Shalon,D., Davis,R.W. and Brown,P.O. (1995) Quantitative monitoring of gene expression patterns with a complementary DNA microarray, *Science*, **270**, 467-470.

Figure 1. Screenshots from the *Machaon CVE* explaining different aspects of its functionality.
(a) Data set window with parameter window for hierarchical clustering.
(b) Result tree illustrating clustering and validation results.