

***Bioinformatics for the management, analysis
and interpretation of microarray data***

Bologna, November 27 - 28, 2003

***Cell classification using gene expression data:
a strategy based on
independent component analysis***

D.G. Calò, G. Galimberti, M. Pillati, C. Viroli

***Dipartimento di Scienze Statistiche, Università di Bologna
{calo,galimberti,pillati,viroli}@stat.unibo.it***



KEY FEATURES OF GENE EXPRESSION DATA

“LARGE p , SMALL n ” PROBLEM:

- *Very large number of variables (genes): from 5,000 to 10,000*
- *Small number of observations (cells): less than 100*

MANY VARIABLES ARE NOISY OR NOT RELEVANT TO CLASS PREDICTION

- *Pearson correlation coefficient*
- *Distance based metrics*
- *Class separation: $BSS(X_i)/WSS(X_i)$*



KEY FEATURES OF GENE EXPRESSION DATA

COMPLEX INTERACTIONS BETWEEN GENES; REDUNDANCY

Genes are points in a n -dimensional space

MOST OF THE ABNORMALITIES IN CELL BEHAVIOUR
ARE DUE TO IRREGULAR GENE ACTIVITIES

Look for “outlying” genes

GENE EXPRESSION PROFILES ARE TYPICALLY NON-GAUSSIAN

Independent Component Analysis



INDEPENDENT COMPONENT ANALYSIS (ICA)

THE MODEL:

$$\mathbf{X} = \mathbf{A}\mathbf{S}$$

$\mathbf{X}=(X_1, X_2, \dots, X_n)$ observed variables ($E(\mathbf{X}) = \mathbf{0}$)

$\mathbf{S}=(S_1, S_2, \dots, S_k)$ latent variables ($E(\mathbf{S}) = \mathbf{0}$)

$\mathbf{A}= n \times k$ mixing matrix, $k \leq n$

THE VARIABLES S_j ARE ASSUMED:

- To be statistically independent:

$$\text{pdf}(S_1, S_2, \dots, S_k) = \prod \text{pdf}(S_j)$$

- To have non-Gaussian distributions

INDEPENDENT COMPONENT ANALYSIS (ICA)

$$Y = WX$$

Restrictions on the extracted variables Y_j :

- $E(Y_h Y_i) = E(Y_h)E(Y_i) \quad h \neq i \quad h, i = 1, \dots, k$
(in order to reduce the number of free parameters)
- $E(Y_j^2) = 1$
(conventional assumption)

For uncorrelated variables:

- $I(Y_1, Y_2, \dots, Y_k) = J(Y_1, Y_2, \dots, Y_k) - \sum J(Y_j)$
(where I denotes mutual information and J denotes negentropy)

Therefore,

the less dependent are the most non-Gaussian ones

ICA leads to meaningful results whenever the probability distribution of X is far from Gaussian

ICA IN GENE EXPRESSION DATA ANALYSIS:

- G. Hori *et al.* (2002) - *ISMB2002*
- W. Liebermeister (2002) - *Bioinformatics*
- X. Liao *et al.* (2002) - *ICASSP 2002, IEEE International Conference*



THE PROPOSED SOLUTION

- k independent components are extracted from the training set
- the p genes are sorted in increasing order according to their absolute scores on each component
- these k marginal rankings are summarized by taking, for each gene, the highest value
- select the m ($m \ll p$) genes located in the last m positions of this joint ranking



AN APPLICATION TO A REAL DATA SET

THE LYMPHOMA DATA SET (Alizadeh et al., 2000):

$C = 3$ classes

$p = 4026$ genes

$n = 62$ cells ($n_1=11, n_2=9, n_3=42$)

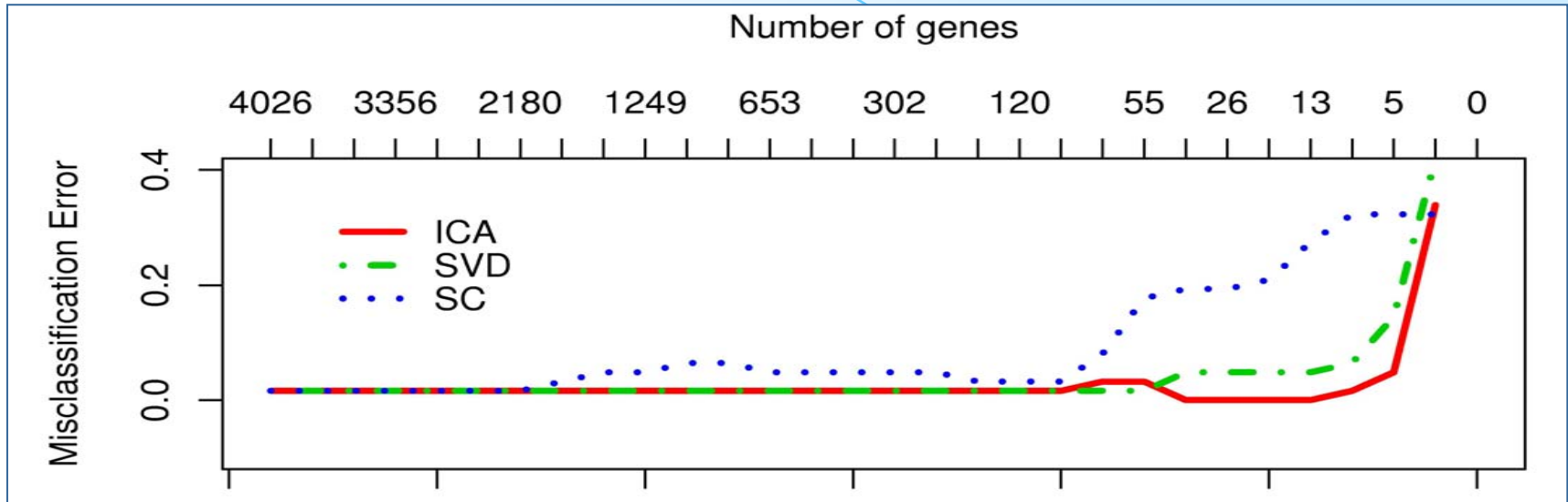
NEAREST SHRUNKEN CENTROIDS METHOD (Tibshirani et al., 2002)

Nearest centroid classification after shrinking each centroid to the overall one.

If a gene is shrunk to zero for all classes, then it is dropped from the allocation rule.

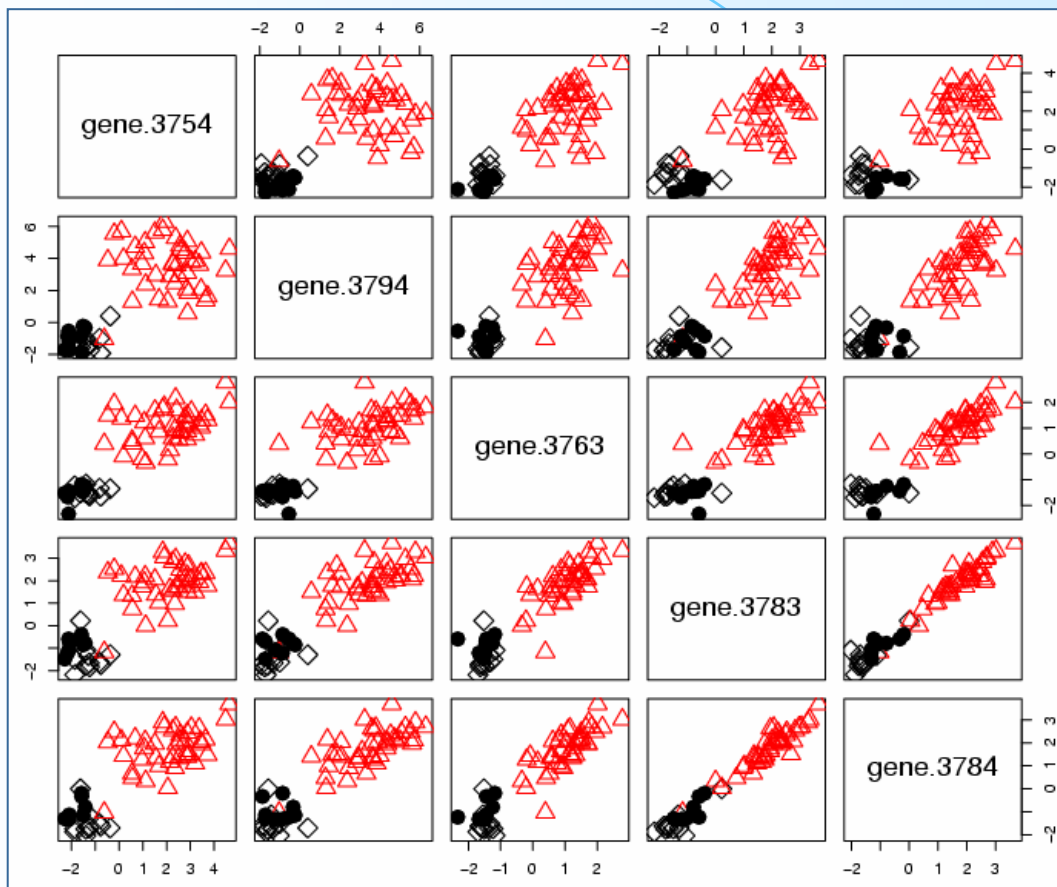


Lymphoma data set: cross-validated misclassification rates (as function of m)

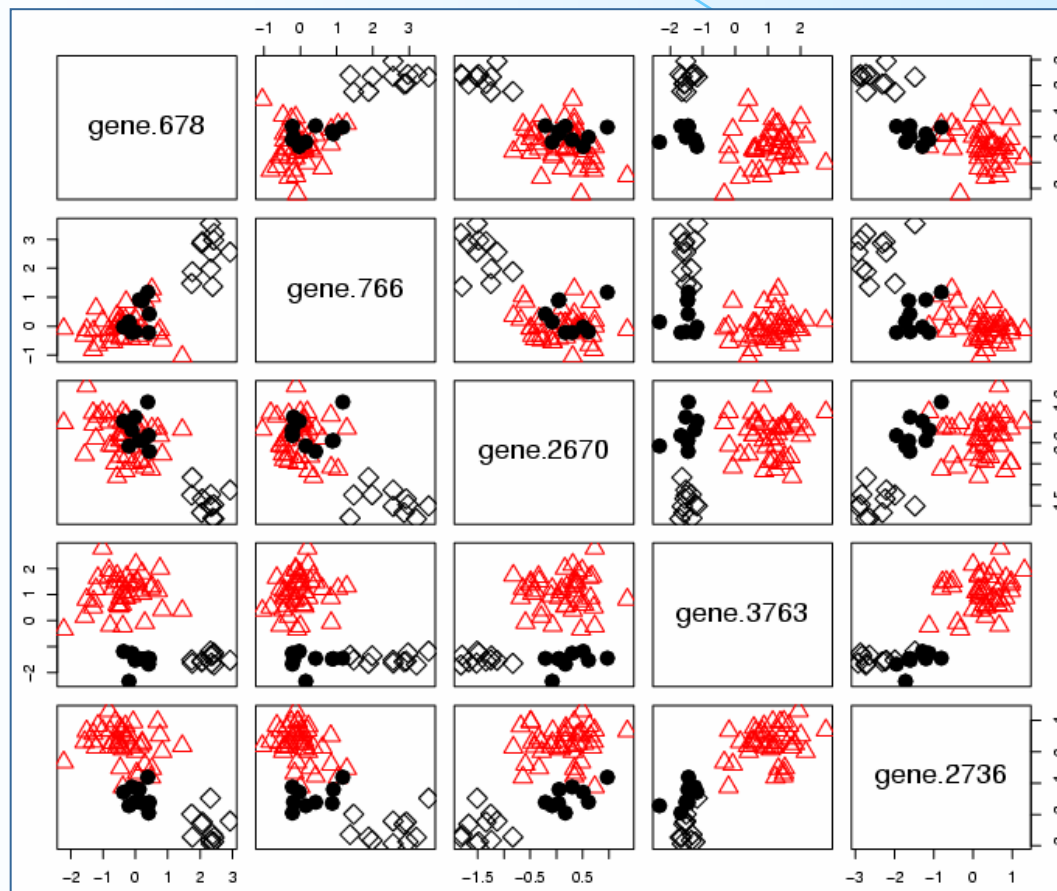


m	222	164	120	93	72	55	34	26	19	13	10	5
SC	0.048	0.032	0.032	0.032	0.081	0.177	0.194	0.194	0.210	0.274	0.323	0.323
SVD	0.016	0.016	0.016	0.016	0.016	0.016	0.048	0.048	0.048	0.048	0.065	0.145
ICA	0.016	0.016	0.016	0.016	0.032	0.032	0.000	0.000	0.000	0.000	0.016	0.048

Lymphoma data set: scatter plot matrix of the last 5 genes surviving the shrinkage procedure



Lymphoma data set: scatter plot matrix of the last 5 genes of the ranking obtained by ICA



OPEN ISSUES

- *Alternatives to the criterion for building the joint ranking*
- *How to choose the number k of the components*
- *How to choose the number m of retained genes*
- *Possible interactions between the proposed selection method and different allocation rules*

