
Cancer profiles by affinity propagation

D. Soria and J.M. Garibaldi

School of Computer Science,
University of Nottingham, Jubilee Campus,
Wollaton Road, Nottingham, NG8 1BB, UK
E-mail: dqs@cs.nott.ac.uk
E-mail: jmg@cs.nott.ac.uk

F. Ambrogi* and P. Boracchi

Istituto di Statistica Medica e Biometria 'GA Maccacaro',
Università degli Studi di Milano,
Via Venezian 1, 20133 Milano, Italy
E-mail: federico.ambrogi@unimi.it
E-mail: Patrizia.boracchi@unimi.it
*Corresponding author

E. Raimondi

Struttura di Statistica e Biometria,
Fondazione IRCCS Istituto Nazionale Tumori,
Via Venezian 1, 20133 Milano, Italy
E-mail: ele.raimondi@gmail.com

E. Biganzoli

Istituto di Statistica Medica e Biometria 'GA Maccacaro',
Università degli Studi di Milano,
Via Venezian 1, 20133 Milano, Italy
and
Struttura di Statistica e Biometria,
Fondazione IRCCS Istituto Nazionale Tumori,
Via Venezian 1, 20133 Milano, Italy
E-mail: elia.biganzoli@unimi.it

Abstract: The affinity propagation algorithm is applied to various problems of breast and cutaneous tumours subtyping using traditional biologic markers. The algorithm provides a procedure to determine the number of profiles to be considered. Well-known breast cancer case series and cutaneous melanoma were used to compare the results of the affinity propagation with the results obtained with standard algorithms and indexes for the optimal choice of the number of clusters. Results from affinity propagation are consistent with the results already obtained having the advantage of providing an indication about the number of clusters.

Keywords: affinity propagation; clustering algorithms; breast cancer; cutaneous melanoma.

Reference to this paper should be made as follows: Soria, D., Garibaldi, J.M., Ambrogi, F., Boracchi, P., Raimondi, E. and Biganzoli, E. (2009) ‘Cancer profiles by affinity propagation’, *Int. J. Knowledge Engineering and Soft Data Paradigms*, Vol. 1, No. 3, pp.195–215.

Biographical notes: Daniele Soria received his BSc and MSc in Applied Mathematics from the University of Milan, Italy, in 2004. He is currently working towards his PhD at the Intelligent Modelling and Analysis (IMA) Research Group, School of Computer Science, University of Nottingham, UK. His current research interests include bioinformatics and medical applications.

Jonathan M. Garibaldi received his BSc (Hons.) in Physics from University of Bristol, UK, in 1984, and MSc and PhD from the University of Plymouth, UK, in 1990 and 1997, respectively. He is currently an Associate Professor and Reader with the Intelligent Modelling and Analysis (IMA) Research Group, School of Computer Science, University of Nottingham, UK. His current research interests include modelling uncertainty in human reasoning and complex data analysis, particularly in medical domains.

Federico Ambrogi is Senior Researcher at the Facoltà di Medicina e Chirurgia, Università degli Studi di Milano and received a degree in Economics in 1997 (Summa Cum Laude) at the Bocconi University of Milano. During 1998–2001, he was Consultant Researcher for data analysis for support decision systems at NUNATAC, SAS Institute quality partner. From 2001 to 2005, he was a Research Fellow at the Istituto Nazionale Tumori, Milan. From 2006, he worked as a Senior Researcher at the Institute of Medical Statistics and Biometry, University of Milan. His main interests regard survival regression models for clinical decision support and multivariate techniques for bioprofile identification of complex diseases.

Patrizia Boracchi graduated at the Faculty of Biology, University of Milan in 1983. She obtained her PhD in Medical Statistics in 1989 at the University of Milan. From 1991 to 2000, she was a Technical Collaborator at the Institute of Medical Statistics and Biometry at the faculty of Medicine and Surgery, University of Milan. She is Researcher at the same Institute since 2001. Her main research fields are statistical methods for survival analysis with particular reference at flexible regression modelling of the hazard function and competing risks. She is the author of more than 90 scientific publications on international journals.

Elena Raimondi received his MSc in Biostatistics and Experimental Statistics (Summa Cum Laude) in 2009. Her final Masters project was on ‘Clustering by affinity propagation: a profiling of breast cancer based on biomarkers’, conducted at the Fondazione IRCCS Istituto Nazionale Tumori, Milan, Italy, under the supervision of Dr. Federico Ambrogi and Prof. Elia Biganzoli. She is currently working as a Biostatistician at the Laboratory of Medical Statistics, Department of Cardiovascular Research, Istituto di Ricerche Farmacologiche Mario Negri, Milan, Italy.

Elia M. Biganzoli graduated with a degree in Biology in 1988 from the University of Milan. He holds a PhD in Medical Statistics (1996). In 1995, he served as Senior Biostatistician at Istituto Nazionale Tumori, Milano. Since 2007, he joined the Faculty of Medicine, University of Milan where he is currently Professor of Biostatistics. His main research fields concern statistical methods for survival analysis and biological assay. He developed statistical approaches for the extension of generalised linear models with artificial neural networks and splines for the flexible analysis of censored survival data. He is author of more than 40 scientific publications on international journals.

1 Introduction

Genomic analysis renewed interest in clustering techniques. After the seminal paper of Eisen et al. (1998), proposing hierarchical clustering and the visual inspection of the dendrogram to discover unknown pattern of gene associations, the use of clustering has become more and more popular especially for discovering profiles in cancer with respect to high-throughput genomic data. Important applications of the Eisen method are the work of Bittner et al. (2000) on clustering of cutaneous melanoma and the works of van't Veer et al. (2002) and Perou et al. (2000) on breast cancer.

Recently a classification of breast carcinoma using traditional tumour markers was proposed (Ambrogi et al., 2006). The classification was in agreement with the classifications obtained with c-DNA microarray data (Perou et al., 2000; van't Veer et al., 2002). Different clustering algorithms were used to choose a stable solution across different clustering methods. At last a classification in four clusters was preferred and suggested a possible separation of high risk profiles. One of the main problems connected with cluster analysis is the choice of the number of clusters. In classical cluster analysis it is customary to use indexes to compare one cluster solutions to other cluster solutions and to choose the one suggested as optimal.

In the previous application (Ambrogi et al., 2006), different indexes were used to select an optimal partition. Namely the indexes proposed by Calinski and Harabasz (1974), Krzanowski and Lai (1985), Hartigan (1975) and Tibshirani et al. (2001), were considered. It is worth noting that the visual inspection of the dendrogram suggested by Eisen et al. (1998) is an informal method to determine the number of clusters. Such a procedure was criticised in Goldstein et al. (2002) as it can cause difficulty in assessing the validity of the grouping.

According to Getz et al. (2000), the number of clusters should be determined internally by the clustering algorithm and should not be externally prescribed.

In this work a new clustering algorithm, the affinity propagation (AP) Frey and Dueck, (2007), will be adopted to cluster cancer patients in order to evaluate its performance with respect to the traditional applications. Several datasets, already published in literature, were considered in our study: the melanoma data of Bittner et al. (2000) and four different breast cancer data analyses: namely, the studies of Ambrogi et al. (2006), Perou et al. (2000), van't Veer et al. (2002) and Sotiriou et al. (2003). Although this algorithm does not determine automatically the number of clusters it provides a consistent method to suggest the number of clusters to be created which can be useful to detect different levels of association pattern.

2 Material and methods

2.1 Case data

This section includes a description of the case series adopted in the work.

The information on 633 patients operated on for primary infiltrating breast cancer between 1983 and 1992, archived at the Pathology department of the University of Ferrara, was retrospectively analysed in the work of Ambrogi et al. (2006).

The available data concerned patient age, pathological tumour size, histologic type, pathologic stage and number of metastatic auxiliary lymph nodes; as well as

immunohistological determinations of oestrogen receptor status (ER), progesterone receptors status (PR), Ki-67/MIB-1 proliferation index (Ki-67), *c-Erb-B2*/NEU (NEU) and the p53 oncosuppressor gene (p53).

Values of ER, PR and NEU tended to be grouped on the following values: 0%, 10%, 25%, 50%, 75% and 100%; they were consequently discretised on those values. Values of Ki-67 and p53 were used as originally measured.

A second dataset was also analysed: the melanoma data of Bittner et al. (2000). These data consist of gene expression profiles obtained on a collection of 38 samples, comprised of 31 melanoma tumours and seven controls. For the analysis described in Section 3, the data from the seven control specimens were excluded and only the ratios for the 3613 genes that were considered 'well measured' (that is their intensities were sufficiently high) were used. These ratios were converted to log₂ ratios.

Another dataset on breast cancer that was considered for investigation was the one of Perou et al. (2000). Variation in gene expression patterns in a set of 65 surgical specimens of human breast tumours from 42 different individuals have been characterised, using complementary DNA representing 8102 human genes. According to the authors, these patterns provided a distinctive molecular portrait of each tumour (Perou et al., 2000). Sets of co-expressed genes were identified for which variation in messenger RNA levels could be related to specific features of physiological variation. The tumours could be classified into subtypes distinguished by differences in their gene expression patterns. In their paper, Perou et al. (2000) focused on a set of 1,753 genes in 65 experimental tissue samples. In each sample, the ratio of the abundance of transcripts of each gene to its median abundance across all tissue samples, is reported.

Data and original analyses of Bittner and Perou are fully described in the book *Design and Analysis of DNA Microarray Investigations* by Simon et al. (2004).

The dataset analysed in van't Veer et al. (2002) was also considered in this study: they used DNA microarray analysis on primary breast tumours of 117 young patients and applied supervised classification to identify a gene expression signature strongly predictive of a short interval to distant metastases in patients without tumour cells in local lymph nodes at diagnosis. All patients were lymph node negative and under 55 years of age at diagnosis. Unsupervised clustering detected two subgroups of breast cancer, which differ in ER status and lymphocytic infiltration.

The last dataset we analysed was taken from Sotiriou et al. (2003) where tumour samples from 99 patients with primary local breast cancer were accessed from the John Radcliffe Hospital from January 1993 to December 1994. All of the tumour samples were invasive ductal carcinomas; 46 individuals were node negative and 53 were node positive. Almost all of the patients received adjuvant treatment after surgery, consisting of radiotherapy (80 patients), chemotherapy (32 patients) and endocrine therapy (78 patients) according to accepted practice guidelines at that time. From each woman belonging to the cohort, a tumour tissue sample, on which all the analyses were run, was extracted using microarray technology. For each spot, Log (base 2) ratios between red signal and green signal were calculated. The log ratios were then normalised within each array by subtracting from each the median log ratio value across the spots on the array.

2.2 *Statistical methods*

The clustering technique AP (Frey and Dueck, 2007) will be adopted for grouping tumours with similar biological characteristics.

As other clustering algorithms, this method uses data to find a set of centres such that the sum of squared errors between data points and their nearest centre is small.

Like other traditional clustering techniques, the AP algorithm determines the centres from real data points (exemplars). These exemplars correspond, for example, to the medoids in the algorithm partitioning around medoids (Kaufman and Rousseeuw, 1990) (PAM, a more robust version of K-means), that is k representative objects among the observations of the dataset that should represent the structure of the data.

As a technical detail, it is worth noting that K-means algorithm does not use exemplars, as the centres are not generally actual data points.

AP combines the properties of different classes of clustering algorithms. On one hand, algorithms like hierarchical clustering are based on grouping pairs of objects with high affinity. On the other hand model-based clustering uses a probability model based on a mixture of class conditional distributions. AP uses both pairs comparison and a probability model to determine the optimal grouping. According to a more technical point of view, AP can be derived as the sum-product algorithm in a graphical model describing the mixture model (Frey and Dueck, 2006).

The first step for the algorithm implementation is to choose a measure of similarity, $s(i,k)$, between all pairs of data points. In AP terminology, $s(i,k)$ quantifies how well the data point with index k is suited to be the exemplar for data point i . Generally, as similarity, it is used the negative Euclidean distance. In the case of c-DNA data the Pearson correlation is generally used as similarity measure (Bittner et al., 2000).

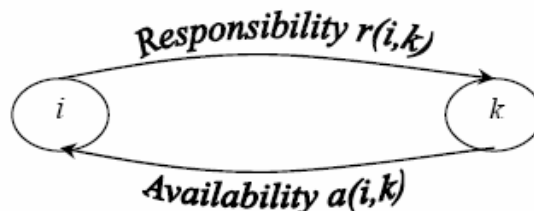
This method does not require the number of clusters to be prespecified.

The second step is about the choice of the values of 'preferences' which will be indicated, with a little abuse of notation as $s(i,i)$. Please note that this is not a similarity measure. The preferences represent a measure of how much data point i is candidate to be an exemplar. In general, data points with larger values of $s(i,i)$ are more likely to be chosen as exemplars. At the beginning, the AP simultaneously considers all data points as potential exemplars (input the preferences common for all data points).

The number of identified exemplars is influenced by the values of the input preferences, but also emerges as a result of the message passing structure that is illustrated subsequently. For very small value of input $s(i,i)$, for every i , all data points are grouped in one large cluster with a single exemplar; in the opposite case of large $s(i,i)$ for every i , each data point prefers to be its own exemplar. In general, the initial value of the preferences is set equal to the median of all input similarities (resulting in a moderate number of clusters) or to their minimum (resulting in a small number of clusters).

The AP is a method that recursively transmits messages (that will be defined subsequently) between pairs of data points until a good set of exemplars and corresponding clusters emerges.

Figure 1 Transmission of messages during affinity propagation



The algorithm is named AP because at any point in time, each message reflects the current affinity between one data point and the other that is its exemplars.

In practice, it is adopted a message-passing algorithm in which each data point i furnishes a measure to suggest another data point k to be selected as cluster centre, taking into account other potential exemplars for point i .

There are two kinds of message being passed between each pairs of data points that represent the relationship between data points:

- *responsibility*: sent from data point i to candidate exemplar k . It is a measure that quantifies how well-suited point k is to be the exemplar for point i , taking into account other potential exemplars for point i . This message is represented by $r(i,k)$ and it is computed using this formula:

$$r(i,k) = s(i,k) - \max_{k' \text{ s.t. } k' \neq k} \{a(i,k') + s(i,k')\},$$

where s.t. means ‘so that’

- *availability*: sent from candidate exemplar point k to point i . It is a measure that reflects the evidence for point i to choose point k as its exemplar, considered that other points may have k as an exemplar. This message is represented by $a(i,k)$ and it is computed using this formula:

$$a(i,k) = \min \left\{ 0, r(k,k) + \sum_{i' \text{ s.t. } i' \notin \{i,k\}} \max \{0, r(i',k)\} \right\}.$$

A particular measure is the ‘self-responsibility’, that is $r(k,k)$; it reflects accumulated evidence that point k is an exemplar and how it would be unsuitable to be integrated in a group of another cluster centre.

At the beginning of the algorithm, the availabilities are initialised to zero, so $r(i,k)$ is set to the input similarity between point i and its potential exemplar k minus the largest of the similarities between point i and other candidate exemplars. After the computation of all the responsibilities, the availabilities are worked out using the previous formula. Only the positive portions of responsibilities between the candidate exemplar k and other data points i' are added because it is only necessary for a good exemplar to explain some data points well ($r(i',k) > 0$) regardless of how poorly it explains other data points ($r(i',k) < 0$). In fact, if $r(i',k) < 0$, k is not suited to be the exemplar for point i' . So in this case, the point i' will not contribute to the message passing from candidate exemplar k to point i .

After that, the messages are recursively updated for a fixed number of iterations or until a stable clustering result. At any stage, the availabilities and responsibilities can be combined to identify exemplars. For point i , the value k that maximises $a(i,k) + r(i,k)$ identifies point i as exemplars if $k = i$ or identifies the data point that is the exemplars for point i . In other words, as suggested in Dueck et al. (2008), after exchanging messages, AP identifies a set of exemplars K so as to maximise the net similarity, which is defined by the authors as

$$\sum_{i \notin K} \max_{k \in K} s(i,k) + \sum_{k \in K} p(k)$$

where $p(k)$ is the *a priori* preference that point k be chosen as an exemplar.

At the end of the message passing, we obtain the number of clusters and the labels for each data point of its exemplars.

All the equations shown above are derived and explained in details in the supporting online material for Frey and Dueck (2007).

An advanced characteristic of AP is that it determines the number of clusters on the basis of the message passing architecture and the points that are most representative, given an initial common preference. It is possible to see the effect of the value of the input preference on the number of clusters by a graphic with the value of the common initial preference on the x-axis and the respective number of clusters on the y-axis. In this way, the value to adopt in the analysis can be established in correspondence with plateaus that are observable in this graphic.

Given the initial common preference AP defines a unique solution. One of the strong points of AP is its computational efficiency, as described in Leone and Weigt (2007). The algorithm is feasible even in presence of very large data sets.

In all our applications of AP, 1-Pearson correlation was used as similarity. The only exception was in the first case series where the negative Euclidean distance was used. These choices were in accordance with the ones adopted by the authors of the works considered.

Multiple correspondence analysis (MCA), see Greenacre (1994) or Lebart et al. (1995), was used as visualisation technique to study the composition of the clusters for the breast cancer data due to the discretisation of the values of the biomarkers. The five biologic markers (ER, PgR, Ki-67, NEU and p53) were used to create the MCA plot (active information). The cluster classifications were used as passive information. The amount of information explained by the first two axes was calculated following the approach suggested by Benzecri (1979). In fact, due to a geometric property of MCA, the percentages of the inertia explained by each axis are always a pessimistic indicator of the quality of the representation. Therefore, Benzecri suggested the following indicator:

$$\varphi(\lambda) = \left(\frac{p}{p-1} \right)^2 \left(\lambda - \frac{1}{p} \right)^2$$

where p is the number of variables and λ is the principal inertia. For the melanoma data principal component plots (Venables and Ripley, 2002) were used to visualise the separation of the tumour samples according to the c-DNA microarray data.

3 Results

3.1 Ambrogi et al. breast cancer biomarkers data

AP was applied to the breast cancer data analysed in Ambrogi et al. (2006) where the final clustering was obtained using a K-medoids algorithm to generate four clusters.

A graphical evaluation of the effect of the value of the input preference on the number of clusters for the breast cancer data is reported in Figure 2.

The presence of three main plateaus in correspondence with two, four and five clusters is shown.

When doing the analysis with two clusters (results not shown), by using an input preference value in correspondence with that plateau, results consistent with expectations

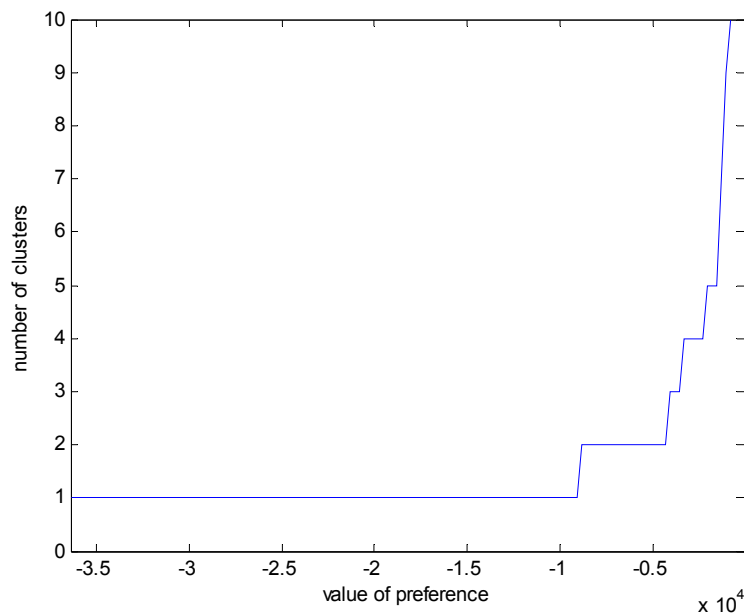
tied to data from literature were obtained. Indeed, one cluster was associated with null values of ER and PR, the other with high values of these biological markers.

Then the message-passing algorithm was run with an input preference to obtain four clusters. The results are reported in the MCA plot (Figure 3). The information explained by the first two axes is near to 89%. Therefore, the two-dimensional plots are expected to be effective representations of the associations displayed.

The MCA plot was generated with the categories of the five biological markers as active information and the AP cluster classification as passive.

As for the contribution of the categories of the biological markers to the construction of the MCA axes, along the first axis there was a separation between high values of PR, ER and the categories of ER and PR absent, high NEU, P53 and Ki-67. The second MCA axis mainly separated the highest values of ER and PR from low PR, Ki-67 and null category of p53.

Figure 2 The effect of the value of the input preference on the number of clusters for the Ambrogi et al. (2006) data (see online version for colours)



Null values of PR and high values of NEU were associated with the Cluster 4. Null values of ER, highest values of p53, Neu and Ki-67 were associated with Cluster 3. Therefore, Cluster 3 and Cluster 4 represent groups that are associated with characteristics known to be poor prognostic factors. Whereas, Cluster 1 was associated with highest values of ER and PR, so it seemed to represent subject with characteristics known to be good prognostic factors. Cluster 2 seemed to be associated with intermediated values of PR and ER and null values of Neu; so also this cluster was associated with less aggressive tumour features. As for the triple negative patients, null values of PR, ER and NEU associated with positive values of p53 were grouped in Cluster 3.

The distribution of subject between the classification using K-medoids and the classification using AP is reported in Table 1.

If these results were compared with those of the previous work, null values of PR, ER, NEU and p53 were grouped in Cluster 2, which was the cluster most similar to the characteristics of total sample. Instead, in this new classification null values of PR, ER and NEU associated with null values of p53 lay in Cluster 4, a cluster that is not similar to total sample for the distribution of biological markers and represents groups with poor prognostic factors.

Afterwards, the AP algorithm was applied again to obtain a division of subjects in five groups and to compare these results with the four clusters. To do this, a preference value from the plateau in correspondence of five clusters in the first graphic was chosen.

The distribution of subjects between the classification using K-medoids and the classification using AP is reported in Table 2.

Cluster 4 was more associated with PR absent and high values of NEU. Cluster 2 seemed to be more associated with intermediated values of PR and ER and null values of Neu; it was more associated also with low values of KI-67.

As before Cluster 1 and Cluster 2 were associated with less aggressive tumour features, whereas Cluster 3 and Cluster 4 represent groups with a negative prognosis.

Cluster 5 was mainly characterised by PR absent and high value of NEU. Null values of PR, ER and NEU associated with positive values of p53 were grouped in Cluster 3. Unlike the previous classification, when we divided subjects in five groups, null values of PR, ER and NEU associated with null values of p53 move from Cluster 4 to Cluster 5.

Table 1 The distribution of subjects between new and old classification

		<i>AP clusters</i>			
		<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>
<i>Previous work's clusters</i>	<i>1</i>	253	1	0	2
	<i>2</i>	1	122	0	84
	<i>3</i>	0	1	88	2
	<i>4</i>	1	25	1	52

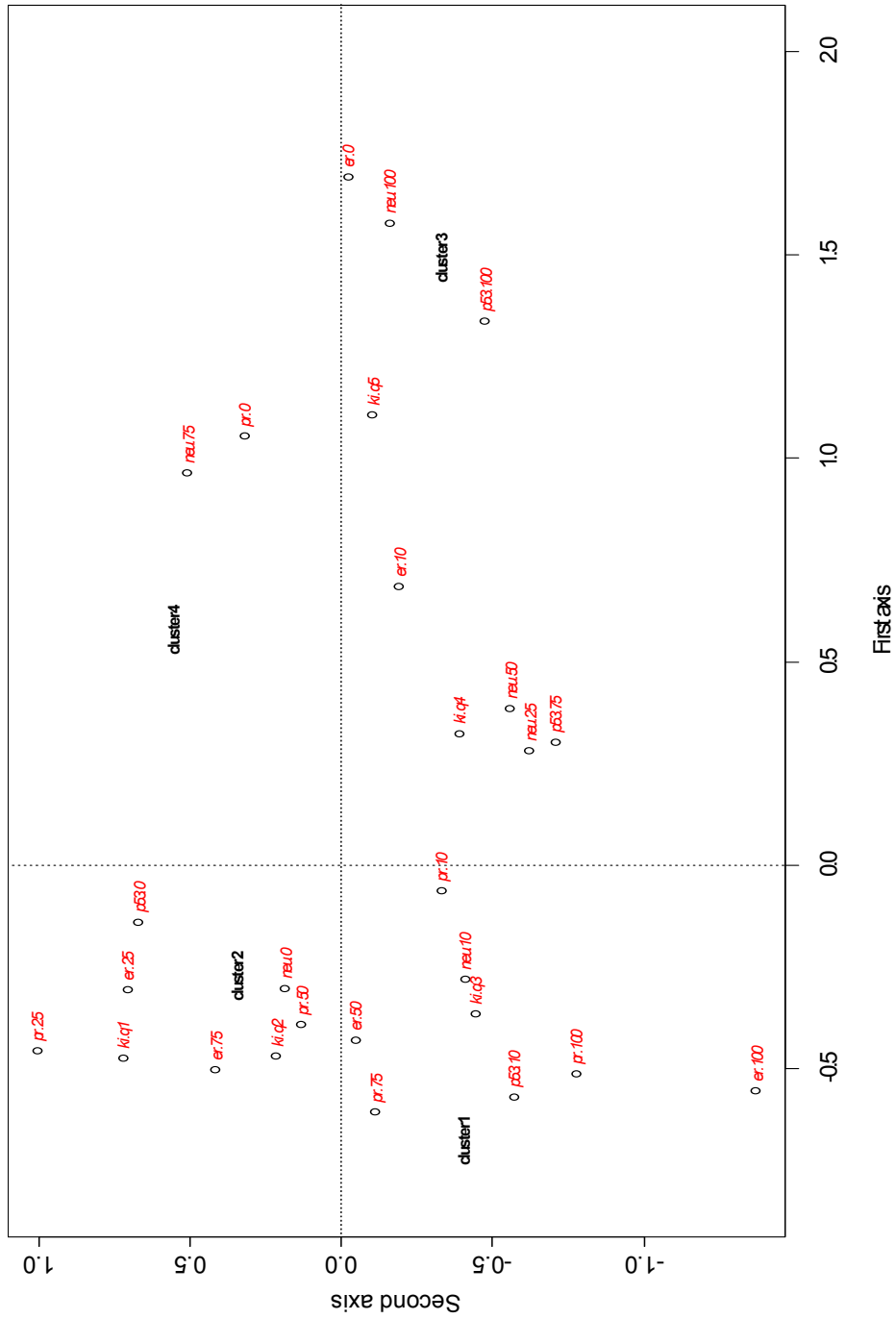
Source: Ambrogi et al. (2006) data

Table 2 The distribution of subjects between new and old classification

		<i>AP clusters</i>				
		<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
<i>Previous work's clusters</i>	<i>1</i>	211	40	0	1	4
	<i>2</i>	0	123	0	0	84
	<i>3</i>	0	1	87	0	3
	<i>4</i>	1	2	0	74	2

Source: Ambrogi et al. (2006) data

Figure 3 MCA plot of the five discretised biological markers ER, PR, MIB, NEU, p53 (active information) and four clusters (passive information) (see online version for colours)

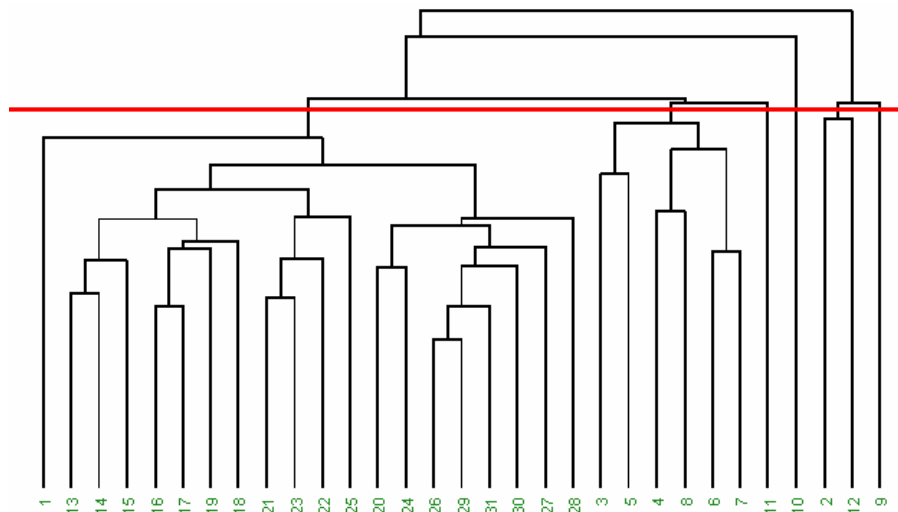


3.2 Bittner et al. melanoma data

Bittner et al. (2000) attempted to determine if c-DNA microarray data could be used to identify distinct subtypes of cutaneous melanoma, a malignant neoplasm of the skin. In particular they were able to identify two major cancer profiles with different biological characteristics. The result was based on the application of a hierarchical algorithm and by cutting the dendrogram by visual inspection (Goldstein et al., 2002).

In Figure 4 the dendrogram resulting from the application of a hierarchical algorithm with average linkage and a similarity matrix based on Pearson correlation is reported. The two clusters were obtained by cutting the tree to obtain five clusters. In this way the 31 melanomas were divided in a single group comprising 20 melanomas while the remaining 11 (actually grouped in four clusters) were considered together.

Figure 4 Dendrogram resulting from the application of hierarchical algorithm to Bittner et al. (2000) dataset (see online version for colours)



AP algorithm was applied to the melanoma data using a distance matrix based on correlations. The resulting plot of the cluster number for different preferences levels is reported in Figure 5. The plot suggests solutions with two, three and five clusters. The solution with five cluster is the one more similar to the one obtained by Bittner et al. (2000). The three-dimensional principal component plot in Figure 6 shows the two groups of the 31 melanomas. The red crosses correspond to the 'interesting' cluster identified by Bittner et al. (2000). The four black squares are tumours classified differently by AP and the hierarchical algorithm. The concordance between the two methods appears satisfying.

Figure 5 The effect of the value of input preference on the number of clusters for the melanoma data (see online version for colours)

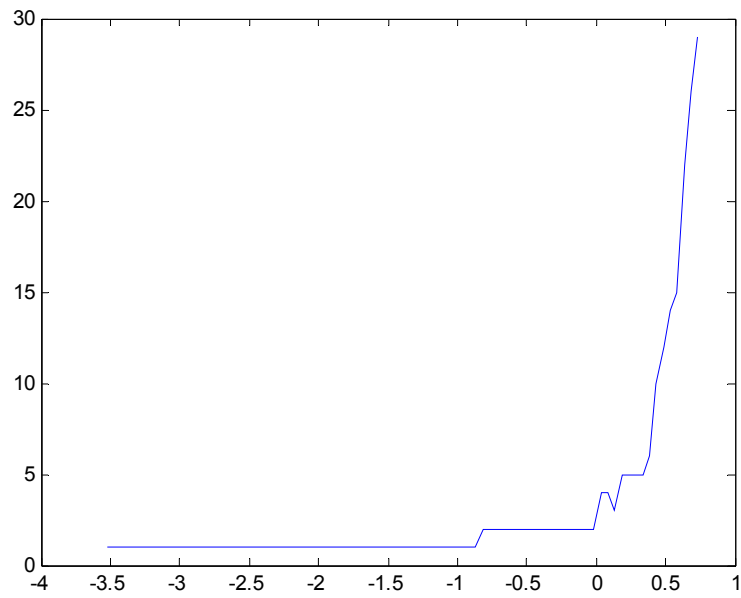
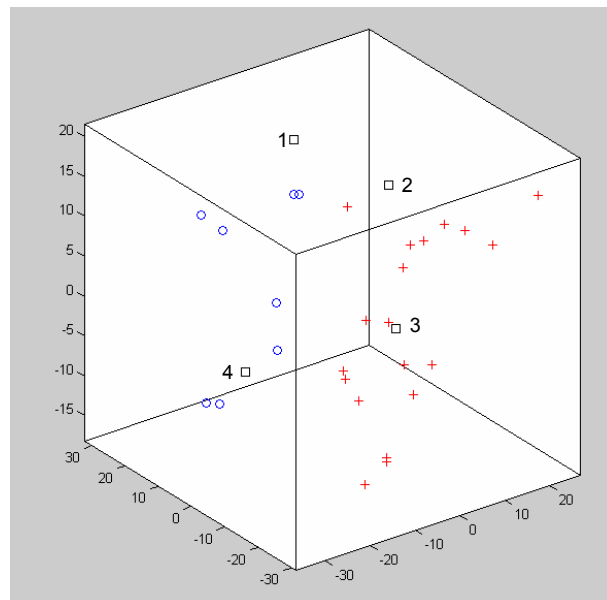


Figure 6 Principal component plot of the gene expression profiles obtained for the 31 melanoma tumours (see online version for colours)



3.3 Perou et al. breast cancer data

In the work of Perou et al. (2000), an average-linkage hierarchical clustering method, as implemented by Eisen et al. (1998), was used to group the experimental samples on the basis of similarity in their patterns of expression. In the dendrogram derived from the hierarchical clustering, two large branches were apparent separating the tumour samples into those that were clinically described as ER positive and those that were ER negative. Within these large branches there were smaller branches for which common biological themes could be inferred, namely basal-like, Erb-B2+, normal-like and luminal epithelial/ER+ cancers.

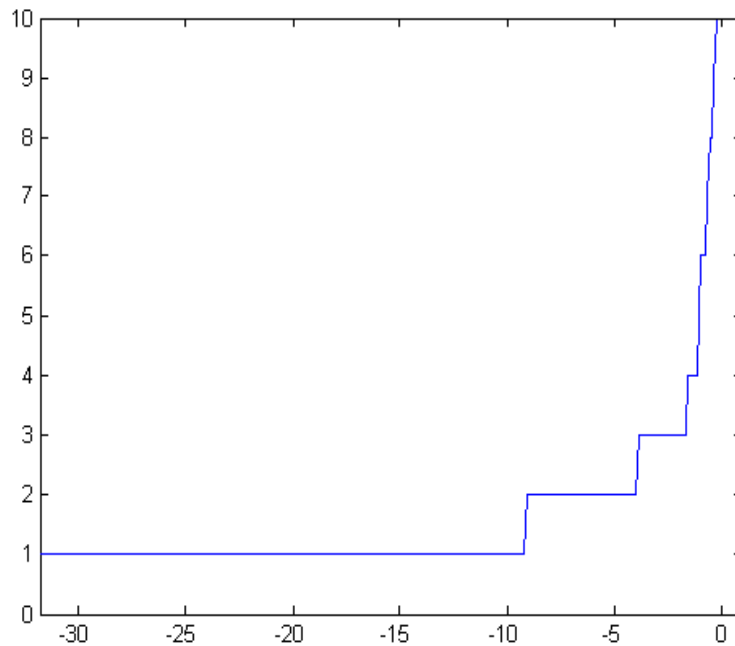
As done in Lama et al. (2004), missing values were imputed with the average values of neighbour genes resulted from a k-nearest neighbours (k-NN) algorithm, setting $k = 10$ (R software package EMV, see Troyanskaya et al., 2001). The ten genes ‘more similar’ to the one with missing value were selected on the basis of the sample correlation: a measure which conforms well to the intuitive biological meaning of ‘co-expression’. In this way, all the unavailable data could be recovered and a 1753×65 matrix is obtained, where rows represented genes and columns represented samples.

The AP algorithm was run and the effect of the value of input preference against the number of clusters is visible in Figure 7.

By the inspection of plateaus, two, three, four and six may be considered as numbers of clusters determining stable solutions.

When considering two clusters a clear correspondence with the ER+ and ER- groups was not evident. In particular the distribution shown in Table 3 was obtained (please note that for five patients information about their ER status was missing).

Figure 7 The effect of the value of input preference on the number of clusters for the Perou et al. (2000) data (see online version for colours)



Considering three clusters, the ER- group was basically identified by the first cluster, which, on the other hand, contained also several ER+ patients (table not shown).

For four clusters a comparison with the four groups determined by the authors was done: from Table 4 it can be seen that all the basal-like tumours were captured by Cluster 1 and the normal-like tumours were assigned to Cluster 2. Clusters 3 and 4 seem to be characterised by the luminal/epithelial cancers. Instead, the Erb-B2 group was not captured by a single cluster by the AP. It is important to note that in their work, Perou et al. (2000) assigned three samples to the ER negative group without specifying which of the four subgroups they are part of.

Finally, when our database was divided in six clusters it was found what is reported in Table 5.

Also with this classification in six groups, the Basal-like patients were assigned to Cluster 1. Instead the Normal ones were divided in Clusters 2 and 5. Clusters 3, 4 and 6 seemed to be three Luminal groups. Moreover, the Erb-B2 group was not captured by a single cluster by the AP.

Table 3 The distribution of subjects between new and old classification

		<i>AP clusters</i>	
		<i>1</i>	<i>2</i>
<i>Previous work's clusters</i>	<i>ER pos</i>	26	22
	<i>ER neg</i>	8	4

Source: Perou et al. (2000) data

Table 4 The distribution of subjects between new and old classification

		<i>AP clusters</i>			
		<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>
<i>Previous work's clusters</i>	<i>Basal</i>	8	0	0	0
	<i>Erb-B2</i>	3	2	1	1
	<i>Normal</i>	0	11	0	0
	<i>Luminal</i>	7	3	16	10

Source: Perou et al. (2000) data

Table 5 The distribution of subjects between new and old classification

		<i>AP clusters</i>					
		<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>
<i>Previous work's clusters</i>	<i>Basal</i>	7	0	0	0	1	0
	<i>Erb-B2</i>	2	2	0	1	0	2
	<i>Normal</i>	0	5	0	0	6	0
	<i>Luminal</i>	2	2	15	10	1	6

Source: Perou et al. (2000) data

3.4 *van't Veer et al. breast cancer data*

In the work of van't Veer et al. (2002), an unsupervised hierarchical clustering algorithm allowed the authors to cluster 117 tumours on the basis of thousands of genes. For clustering, the 'one minus correlation' distance was used by the authors as described in the supplementary information for van't Veer et al. (2002). In this way two subgroups of breast cancer were detected, differing in ER status (ER-positive and ER-negative) and lymphocytic infiltration (see Figure 1a in van't Veer et al., 2002).

Two hundred ninety-three genes had missing information for all 117 patients. These genes were excluded for the analysis. Other missing values were detected in the remaining data and as done for the Perou et al. (2000) data, they were imputed with k-NN algorithm. In this way, our data matrix had 117 tumours and 24,188 genes of log ratio of the intensities of the red and green channels.

Over this dataset, AP algorithm was applied in order to verify the two ER groups, obtaining the plot in Figure 8 showing the effect of the value of the common input preference on the number of clusters.

From the plot above, looking at the plateaus, solutions with two, three, five, six and seven clusters were suggested. It is worth noting that AP was not able to indicate a solution with four clusters but one in six groups was evident in accordance with results previously obtained by Perou et al. (2000).

When considering two clusters a very good concordance with the two groups obtained separating the ER positive and ER negative patients could be seen. In fact, just 16 over 117 cases did not show such a correspondence (see Table 6).

The split in three clusters is nothing more than a subdivision of the previously determined Cluster 1 into two subgroups (table not shown).

Figure 8 The effect of the value of input preference on the number of clusters for the van't Veer et al. (2002) data (see online version for colours)

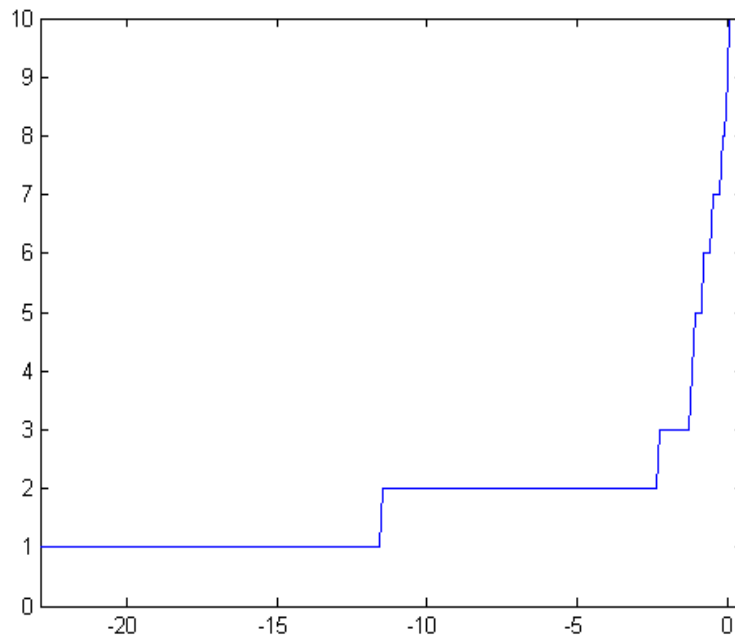


Table 6 The distribution of subjects between new and old classification

		<i>AP clusters</i>	
		<i>1</i>	<i>2</i>
<i>Previous work's clusters</i>	<i>ER pos</i>	64	11
	<i>ER neg</i>	5	37

Source: van't Veer et al. (2002) data

3.5 *Sotiriou et al. breast cancer data*

Cluster analyses were conducted to search for natural groupings in the profiles. As in the original analysis of Sotiriou et al. (2003), before clustering, a screening procedure was applied to eliminate genes showing minimal variation across the set of 99 specimens. Specifically, for each gene, the 5th and 95th percentiles of the ratios were calculated. If the ratio of the 95th to 5th percentile was <3 that gene was not included in the cluster analysis. This process left 706 probe elements for the cluster analysis (Sotiriou et al., 2003).

Authors applied hierarchical agglomerative clustering to these normalised log ratios by using both compact linkage and average linkage and both Euclidean and one minus Pearson correlation distance metrics. Normalised log ratios were median-centred within each gene for all of the cluster analyses. The clustering results obtained by using compact linkage with one minus Pearson correlation distance applied to the 706 probe elements appeared by visual inspection to yield the most distinctive clusters (remaining blinded to any clinical or outcome variables) and hence this was the clustering algorithm used for the unsupervised cluster analyses based on these probe elements.

Hierarchical clustering results showed that the tumour samples could be confidently separated into two main groups primarily associated with ER status as determined by the ligand-binding assay and confirmed by immunohistochemistry. Eighty-two percent of the tumours in the left main branch are ER- and 88% of the tumours in the right main branch are ER+ by immunohistochemistry. This finding corroborated the earlier analysis that the major factor discriminating the expression phenotype is ER status. As expected, the ER- cluster had a higher percentage of high-grade tumours than the ER+ cluster. The dendrogram further branched into smaller subgroups within the ER+ and ER- classes.

In the ER- group, tumours characterised by expressions of genes typical of the 'basal-like' groups were evident. In particular, based on levels of expressions of specific genes, it was possible to distinguish the 'basal 1' group from the 'basal 2' one. Moreover, in the ER- group, a subgroup characterised by the over-expression of the HER2/NEU oncogene was evident in the basal subgroup.

Concerning the ER positive group, subjects in this cohort were characterised by the expressions of genes specific of the 'luminal-like' tumours. This group was further segregated into three smaller subclasses: Luminals 1, 2 and 3.

When considering this dataset for the AP analysis, the same technique previously described for dealing with missing values was used.

As for the other datasets considered, the effect of the value of input preference on the number of clusters is reported in Figure 9.

As it can be seen from the plot, there are three main plateaus in correspondence of solutions with two, three and four clusters and other two smaller ones suggesting six and

seven groups. When considering two groups, results consistent with the ones from literature and the two groups of Sotiriou et al. (2003) are derived. In fact, classification of subjects affected by breast cancer is associated to the ER status: basing on the immunohistochemistry, the first cluster is characterised by 64% of ER negative, while the second by the 92% of ER positive patients (see Table 7). The distribution of subjects in the two clusters with respect to the histological grade is also reported in Table 7.

Table 7 The distribution of subjects between AP clusters, grade and ER status

		<i>AP clusters</i>	
		<i>1</i>	<i>2</i>
<i>Previous work's clusters</i>	<i>Grade 1</i>	2	14
	<i>Grade 2</i>	14	24
	<i>Grade 3</i>	31	14
	<i>ER pos</i>	17	48
	<i>ER neg</i>	30	4

Source: Sotiriou et al. (2003) data

From the table above it can be seen that the ER– group (Cluster 1) is basically formed by those patients with histological grade equal to two or three.

After that, in order to obtain a classification of 99 patients in three groups, a preference value in correspondence of the proper plateau was chosen. With this value as an input, the AP algorithm was run again. From the classification in three clusters, a particular group, Cluster 3, is visible: it is mainly formed by patients with negative ER and high histological grade (tables not reported).

The last significant plateau visible in Figure 9 suggests to consider four groups (tables not shown).

To interpret the four groups obtained with AP on the basis of the genic expression, boxplots of the principal tumour markers were used (plots not shown). This was the only subdivision for which such an analysis could be performed.

Cluster 1 can be characterised as a ‘luminal-group’; this fact is also confirmed by the overexpression of the ER in this group with respect to the overall population behaviour.

Samples belonging to Cluster 2 present an overexpression of the basal cytokeratins (CKs) and a low expression of the luminal ones. Moreover this group is the only one which has low levels of ER. For these reasons Cluster 2 may be interpreted as a basal-like group even though all Erb-B2 samples were classified in this cluster.

Cluster 3 is characterised by an overexpression of CK8, CK19 (luminal markers) and ER. It can be then considered as a ‘luminal-like’ group.

The values distribution of CKs in Cluster 4 is not different from the one in the overall population. This group is also characterised by the positive estrogen receptor status. It may be concluded that AP is not able to distinguish Cluster 4 from the general population in terms of gene expression.

As in previous examples, a small plateau at six clusters is visible in Figure 9. The distribution of patients in six clusters from AP compared with the one from the Sotiriou’s et al. (2003) groups is reported in Table 8.

Sotiriou’s basal groups (basal 1 and basal 2) were both captured by AP Cluster 2, while the Erb-B2 was now associated with Cluster 6. The Luminal groups were assigned to different clusters: in particular, luminal 1 was described by samples in Clusters 4 and

5, luminal 2 in Cluster 3 and luminal 3 was mainly associated to Cluster 1 but presents several cases in other clusters as well.

Figure 9 The effect of the value of input preference on the number of clusters for the Sotiriou et al. (2003) data (see online version for colours)

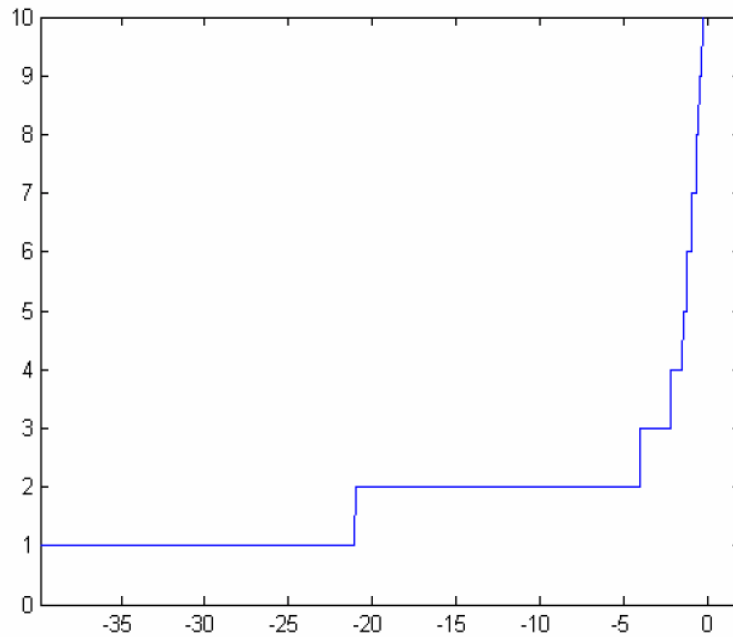


Table 8 The distribution of subjects between new and old classification

		<i>AP clusters</i>					
		<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>6</i>
<i>Previous work's clusters</i>	<i>Basal 1</i>	0	14	0	0	0	2
	<i>Basal 2</i>	0	8	0	1	1	0
	<i>Erb-B2</i>	0	0	0	0	0	7
	<i>Luminal 1</i>	0	0	1	8	10	0
	<i>Luminal 2</i>	9	0	10	2	2	0
	<i>Luminal 3</i>	10	2	3	1	2	6

Source: Sotiriou et al. (2003) data

4 Discussion

Cluster analysis is a powerful technique to explore complex diseases and improve prognosis. The recent literature on omic data is rich of new methods of cluster analysis

able to deal with huge datasets. Moreover techniques of visualisations are usually adopted to suggest the number of clusters (Eisen et al., 1998).

At the same time many papers warn against the possible misuse of clustering techniques (Goldstein et al., 2002).

One of the main problems is the subjectivity of the analysis and the ability of clustering algorithms to create clusters even in absence of real structure.

The choice of the number of clusters is one of the main problems to be faced when applying this kind of analysis. The possibility to use algorithms that incorporate a criterion for the choice of the optimal partition is one of the achievements of the recent developments in this research field. The AP algorithm is characterised by a simple software implementation and it has the ability to suggest the cluster number. In addition, this algorithm has the advantage of taking into account real data points as exemplars; in this way, also, categorical data may be analysed using, as similarity, a different distance from the Euclidean one. By looking for a set of exemplars and not for a set of centroids, the search space is restricted, improving the computational efficiency of the AP. In this work it was demonstrated how the algorithm is in agreement with the solutions obtained with much more effort with traditional algorithms and indexes for the cluster number choice. Moreover the range of the suggested solutions gives insights in the hierarchical structure of the data highlighting different level of information for the treatment of cancer patients well in accordance with previous knowledge. In particular the solution with two clusters for Ferrara breast cancer data, evidenced in Figure 2, reflects the well-known separation between tumours ER positive and negatives. This is a very important distinction and, in fact, in a number of paper of the pre-genomic era the number of clusters considered was in fact two (Querzoli et al., 1995; Ménard et al., 1999). The solution with four clusters is in agreement with the solution selected in the previous work and the four clusters obtained are similar to that created by the PAM algorithm. The solution with five clusters suggests a possible more complex pattern to be explored. The clustering obtained by AP on the melanoma data is able to reproduce the interesting findings of Bittner et al. (2000) having the advantage of avoiding any arbitrary choice due to the visual inspection of the dendrogram. Over Perou et al. (2000) the AP could almost reproduce previous classification: the basal-like tumours were captured by Cluster 1 and the normal-like ones were assigned to Cluster 2; Clusters 3 and 4 represent the luminal-epithelial tumours, while the Erb-B2 group was not identified by AP in a single cluster. The solution with two clusters for van't Veer et al. (2002) data reflects with a few number of exceptions, the separation between ER+ and ER- groups evidenced in the original work. Applying AP over the dataset of Sotiriou et al. (2003), four groups were evident and it was also possible to characterise them on the basis of gene expression of the tumour markers. AP Clusters 1 and 3 were associated with high values of luminal CKs and over expression of ER, while Cluster 2 presents an overexpression of the basal CKs and low levels of ER. For Cluster 4 was not possible to obtain a distinctive characterisation, as it presented similar features to the overall population.

In conclusion, it is worth noting that, for the last three datasets analysed, a common solution in six clusters emerged by the AP algorithm, while one in four groups was not evident in van't Veer et al.(2002) work.

As a future work, the characterisation of the AP groups on the basis of gene expressions will be carried out and compare with the one obtained in the original papers.

Acknowledgements

This study was, in part, funded by the BIOPTRAIN FP6 Marie-Curie EST Fellowship (FP6-007597).

References

- Ambrogio, F., Biganzoli, E., Querzoli, P., Ferretti, S., Boracchi, P., Alberti, S., Marubini, E. and Nenci, I. (2006) 'Molecular subtyping of breast cancer from traditional tumor marker profiles using parallel clustering methods', *Clin Cancer Res*, 1 Feb, Vol. 12, No. 3, Pt 1, pp.781–90.
- Benzécri, J.P. (1979) *Cahiers de l'Analyse des Données*, Vol. 4, pp.377–378.
- Bittner, M., Meltzer, P. and Chen, Y. et al. (2000) 'Molecular classification of cutaneous malignant melanoma by gene expression profiling', *Nature*, Vol. 406, pp.536–540.
- Calinski, R.B. and Harabasz, J. (1974) *Communs Statist*, Vol. 3, pp.1–27.
- Dueck, D., Frey, B., Jojic, N., Jojic, V., Giaeffer, G., Emili, A., Musso, G. and Hegele, R. (2008) 'Constructing treatment portfolios using affinity propagation', Vol. 4955 of *Lecture Notes in Computer Science*, pp.360–371, Springer Berlin/Heidelberg.
- Eisen, M.B., Spellman, P.T., Brown, P.O. and Botstein, D. (1998) 'Cluster analysis and display of genome-wide expression patterns', *Proc Natl Acad Sci USA*, Vol. 95, pp.14863–14868.
- Frey, B.J. and Dueck, D. (2006) 'Mixture modeling by affinity propagation', available at http://books.nips.cc/papers/files/nips18/NIPS2005_0799.pdf.
- Frey, B.J. and Dueck, D. (2007) 'Clustering by passing messages between data points', *Science*, Vol. 315, No. 5814, pp.972–976.
- Getz, G., Levine, E. and Domany, E. (2000) 'Coupled two-way clustering analysis of gene microarray data', *Proc Natl Acad Sci USA*, 24 Oct, Vol. 97, No. 22, pp.12079–12084.
- Goldstein, D.R., Debashis, G. and Conlon, E.M. (2002) 'Statistical issues in the clustering of gene expression data', *Statistica Sinica*, Vol. 12, pp.219–240.
- Greenacre, M.J. (1994) *Theory and Applications of Correspondence Analysis*, Academic Press.
- Hartigan, J. (1975) *Clustering Algorithms*, Wiley, New York.
- Kaufman, L. and Rousseeuw, P. (1990) *Finding Groups in Data*, Wiley, New York.
- Krzanowski, W.J. and Lai, Y.T. (1985) *Biometrics*, Vol. 44, pp.23–34.
- Lama, N., Ambrogio, F., Antolini, L., Boracchi, P. and Biganzoli, E. (2004) 'Some issues and perspective in microarray data analysis in breast cancer: the need for an integrated research', *Proceedings of the 1st European Workshop on the Assessment of Diagnostic Performance*.
- Lebart, L., Morineau, A. and Piron, M. (1995) 'Statistique exploratoire multidimensionnelle', Dunod, Paris.
- Leone, M.S. and Weigt, M. (2007) 'Clustering by soft-constraint affinity propagation: applications to gene expression data', *Bioinformatics*, Vol. 23, pp.2708–2715.
- Ménard, S., Casalini, P., Tomasic, G., Pilotti, S., Cascinelli, N., Bufalino, R., Perrone, F., Longhi, C., Rilke, F. and Colnaghi, M.I. (1999) 'Pathobiologic identification of two distinct breast carcinoma subsets with diverging clinical behaviors', *Breast Cancer Res Treat*, May, Vol. 55, No. 2, pp.169–177.
- Perou, C.M., Sørlie, T., Eisen, M.B., van de Rijn, M., Jeffrey, S., Rees, C.A., Pollack, J.R., Ross, D.T., Johnsen, H., Akslen, L.A., Fluge, Ø., Pergamenschikov, A., Williams, C., Zhu, S.X., Lønning, P.E., Børresen-Dale, A., Brown, P.O. and Botstein, D. (2000) *Nature*, Vol. 406, pp.747–752.
- Querzoli, P., Ferretti, S., Albonico, G., Magri, E., Scapoli, D., Indelli, M. and Nenci, I. (1995) 'Application of quantitative analysis to biologic profile evaluation in breast cancer', *Cancer*, 15 Dec, Vol. 76, No. 12, pp.2510–2517.

- Simon, R.M., Korn, E.L., McShane, L.M., Radmacher, M.D., Wright, G.W. and Zhao, Y. (2004) *Design and Analysis of DNA Microarray Investigations*, Springer.
- Sotiriou, C., Neo, S-Y., McShane, L., Korn, E., Long, P., Jazaeri, A., Martiat, P., Fox, S., Harris, A. and Liu, E. (2003) 'Breast cancer classification and prognosis based on gene expression profiles from a population-based study', *Proc Natl Acad Sci USA*, Vol. 100, pp.10393–10398.
- Tibshirani, R., Walther, G. and Hastie, T. (2001) *J.R. Statist. Soc. B*, Vol. 63, pp.411–423.
- Troyanskaya, O., Cantor, M., Sherlock, G., Brown, P., Hastie, T., Tibshirani, R., Botstein, D. and Altman, R. (2001) 'Missing value estimation methods for DNA microarrays', *Bioinformatics*, Vol. 17, No. 6, pp.520–525.
- van't Veer, L.J., Dai, H., van de Vijver, M.J., He, Y.D., Hart, A.A., Mao, M., Peterse, H.L., van der Kooy, K., Marton, M.J., Witteveen, A.T., Schreiber, G.J., Kerkhoven, R.M., Roberts, C., Linsley, P.S., Bernards, R. and Friend, S.H. (2002) *Nature*, Vol. 415, pp.530–536.
- Venables, W.N. and Ripley, B.D. (2002) *Modern Applied Statistics with S*, 4th ed., Springer.