

# Using Rank-1 Biclusters to Classify Microarray Data

Nasimeh Asgarian, Russell Greiner

Department of Computing Science, and the Alberta Ingenuity Center for Machine Learning,  
University of Alberta, Edmonton, AB, Canada, T6G2E8

## ABSTRACT

**Motivation:** A DNA-microarray measures the gene expression levels of tens of thousands of genes for a particular sample, corresponding to some specific experimental condition. Our goal is to learn a microarray classifier that can distinguish different classes — e.g., to predict which patient will respond well to a treatment, based on the data from his/her microarray. Unfortunately, the large number of genes and the small number of samples make building such classifiers very challenging.

**Results:** This paper proposes a method for learning a microarray classifier by first reducing the dimensionality of the data matrix using *biclusters*, where each bicluster is a subset of genes and a subset of samples whose expression values have similar patterns. We propose a novel algorithm for finding biclusters from the microarray data, based on the best rank-1 matrix approximation, then show how to use these biclusters to classify novel samples. We demonstrate that our method works effectively by comparing its prediction accuracy with that of other classifiers, including one using another bicluster algorithm.

**Contact:** {nasimeh, greiner}@cs.ualberta.ca

## 1 INTRODUCTION

Biologists use DNA-microarrays to measure the gene expression levels of biological samples, where different samples may correspond to different experimental conditions, different states of a tissue (e.g., diseased vs healthy), different individuals, etc. Users analyze such data to understand how expression levels differ in different conditions, which can help determine which genes are involved in some disease, suggest biomarkers of a specific disease, propose targets for drug intervention, and enable the use of microarrays as a screening tool. Such results also have very important applications in pharmaceutical and clinical research (McLachlan *et al.*, 2004).

This paper focuses on classification: e.g., given a patient’s microarray gene expression data, predict whether that patient has cancer or not, or whether that patient will respond to a certain type of treatment or not? Learning such a classifier is challenging for several reasons, including the dimension of the training data: there are typically very few samples with too many features, often less than a hundred samples versus more than 50,000 genes.

One approach around this involves reducing the dimensionality of the data, using perhaps LDA (Somorjai *et al.*, 2003) or SVD (Kluger *et al.*, 2003). We use an alternative, novel way to reduce the dimensionality, based on first finding the “*biclusters*” within the data. Notice this differs from traditional clustering methods, such as hierarchical and *K*-means approaches, that seek various sets of genes that are similar over *all* samples (or perhaps sets of samples that are similar over *all* genes). By contrast, a bicluster is a *specified*

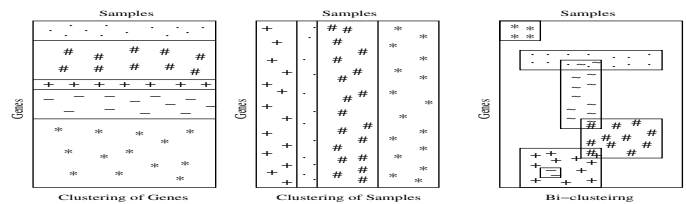


Fig. 1. Clustering versus Biclustering

*subset of genes that are similar for a specified subset of samples*; see Figure 1.

This has a biological motivation, in that we expect that some condition, active in some patients, will trigger some subset of genes, which means that that subset of genes will be co-regulated and co-expressed (that is, correlated) for the patients with that condition. But only for those patients; we expect those genes to behave essentially independently under other conditions (Madeira and Oliveira, 2004).

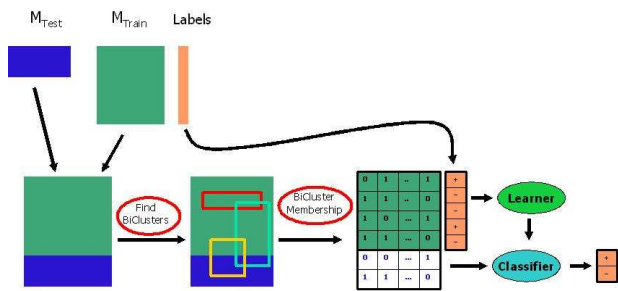
We describe below our particular approach, ROBIC, for finding these correlated genes and patients, based on the *component values* of the principle eigenvectors. There are many other techniques that also use eigenvectors to reduce the dimensionality of a space. These other approaches, however, work by first projecting the instances onto the space defined by the eigenvectors associated with largest eigenvalues; by contrast, our approach explicitly considers the specific *components* of some eigenvectors, and uses these real values to find the (discrete) subset of genes and subset of patients belonging to each bicluster; see Section 3.1.

We then identify each patient with his/her membership in these biclusters, in that bit #*i* is 1 for a patient if this patient belongs in bicluster #*i*. Hence, this process reduces the dimensionality, by mapping each patient-record of say 50,000 real values, to a small number (say 30) of bits. Our system then uses some learning algorithm (e.g., NaïveBayes or linear SVM) to produce a classifier over these 30-tuples; see Figure 2.

Notice this method is implicitly selecting a subset of genes from the original space. For example, if biclusters #5 and #7, based respectively only on genes {*g*<sub>1</sub>, *g*<sub>6</sub>, *g*<sub>93</sub>} and {*g*<sub>1</sub>, *g*<sub>20</sub>, *g*<sub>21</sub>, *g*<sub>45</sub>}, are sufficient to produce the classification, then we are implicitly asserting that these 7 genes are responsible for the class label. (But this might not be based on the actual gene expression values for these genes; see Section 4.2.)

This paper’s two main contributions are...

- An algorithm, ROBIC, for finding biclusters from microarray data, based on the best rank-1 matrix approximation.
- A method, BIC, for using such a biclustering algorithm to learn a sample classifier.



**Fig. 2.** The BIC “Bicluster Classifier” system first finds the biclusters within the training data (without class labels) and test data, then learns a classifier on these “bicluster member” features, which it uses to predict labels for the test data

We demonstrate that this system works effectively across a number of standard DNA microarray classification tasks. The literature shows that many of the “diagnostic” tasks — e.g., predicting whether a patient has some type of cancer — are relatively easy, as a small number of genes are sufficient to produce very high classification accuracy. We focus instead on “prognostic” predictors — e.g., predicting whether a patient (who is known to have cancer) will respond to some specified treatment. These tasks are more difficult both because there tend to be no small set of (say 1 to 5) genes whose expression levels distinguish the classes, and also because the sample sizes are yet smaller, as they include only the subset of patients with the disease.

Section 2 provides background information on biclustering methods and classification tasks, and summarizes previous work related to biclustering and classifying microarrays. Section 3 describes the biclustering problem and the structure of our algorithms: ROBIC for finding biclusters and BIC for using biclusters to classify microarray samples. Section 4 describes the results of applying our biclustering and classification method on a number of publicly available microarray datasets. We also compare our results to those of other approaches (including ones based on an alternative biclustering algorithms) on these datasets, and provide a statistical test (“permutation test”) to demonstrate that our results are significant.

Greiner (2007) provides additional information related to this work, including the datasets used in the studies reported here, as well as other studies, and also other relevant statistics about our results.

## 2 RELATED WORK

### 2.1 Biclustering Algorithms

**2.1.1 The Problem** The microarray data is stored in an  $n \times p$  matrix  $M$ , which is defined by a set of rows (genes),  $G = \{g_1, g_2, \dots, g_n\}$ , and a set of columns (samples),  $S = \{s_1, s_2, \dots, s_p\}$ , where  $|G| = n \gg p = |S|$ . Each element  $m_{ij}$  is a real value, which usually represents (the logarithm of) the expression level of gene  $i$  for sample  $j$ .

A *bicluster* is a subset of genes that have a similar pattern over a subset of samples. This corresponds to a sub-matrix  $M_{I,J}$  of  $M$ , where  $I \subseteq G$  is a subset of genes (rows), and  $J \subseteq S$  is a subset of samples (columns). Our challenge is: Given a data matrix,  $M$ , find (and then use) a set of  $K$  biclusters,  $B_k \approx M_{I_k, J_k}$  where for each

$k = 1, \dots, K$ , the set of  $B_k$ ’s genes  $I_k \subseteq G$  have a similar pattern within the set of  $B_k$ ’s samples  $J_k \subseteq S$ .

**2.1.2 Previous Work** The *Plaid* system (Lazzeroni and Owen, 2002) models gene expression data  $M = (m_{ij})$  using a set of biclusters (which they call “layers”)

$$m_{ij} = \sum_{k=1}^K \theta_{ijk} [\mu^k + \alpha_j^k + \beta_i^k] \quad (1)$$

where  $K$  is number of biclusters,  $\theta_{ijk}$  is 1 if bicluster  $k$  includes position  $(i, j)$  and 0 otherwise, and the values  $\langle \mu^k, \{\alpha_j^k\}_j, \{\beta_i^k\}_i \rangle$  are associated with bicluster  $k$ .<sup>1</sup> Hence, if a bicluster  $k$  includes  $(i, j)$ , then it contributes a value there that is the sum of a background value  $\mu^k$ , a value  $\beta_i^k$  associated with the column  $i$  and a value  $\alpha_j^k$  associated with the row  $j$ . Note that a position can belong to 0, 1 or more biclusters.

While the *Plaid* model is “additive”, we consider a multiplicative variant:

$$m_{ij} = \sum_{k=1}^K \theta_{ijk} [\alpha_j^k \times \beta_i^k] \quad (2)$$

and also allow each position to belong to 0, 1 or more biclusters.

Kluger *et al.* (2003) also seek a multiplicative model. That work, however, considers only the challenge of identifying a single bicluster, but not finding multiple biclusters, nor the interactions between them. Their approach uses *singular value decomposition*, SVD, to obtain a set of eigenvalues and then project patient-tuples onto the associated eigenspaces. While our ROBIC also uses SVD, it does not deal with eigenvalues but instead deals with the *components of the leading eigenvector*. (See Section 2.2 below.)

Neither of these approaches connect their biclusters to a classification task; they instead evaluate the quality of a bicluster by, for example, testing against a null hypothesis of no structure in the data matrix.

Madeira and Oliveira (2004) provide a comprehensive survey that analyzes a large number of existing approaches to biclustering, organizing them according to the type of biclusters they can find, the patterns of biclusters that are discovered, the methods used to perform the search, and the target applications. Also, Preli *et al.* (2006) introduce their new algorithm (Bimax), then provide a systematic comparison and evaluation of prominent biclustering methods, as well as clustering and hierarchical clustering, for the purpose of *gene classification*. They conclude the bicluster-based approaches are the most effective.

### 2.2 Relation to SVD-Based Methods

There are many systems that use SVD to reduce the dimensionality of microarray data (Alter *et al.*, 2000; Wall *et al.*, 2003; Ding, 2003; Hastie *et al.*, 2000). These SVD-based systems first compute a list of eigenvalue/(row)eigenvector/(column)eigenvector triples  $\langle v^{(i)}, \alpha^{(i)}, \beta^{(i)} \rangle_i$  from  $M$ , where here  $v^{(i)} \in \mathbb{R}$  (sorted by  $|v^{(1)}| \geq |v^{(2)}| \geq \dots$ ),  $\alpha^{(i)} \in \mathbb{R}^n$  and  $\beta^{(i)} \in \mathbb{R}^p$ , such that  $M \approx \sum_i v^{(i)} \alpha^{(i)} \times \beta^{(i)T}$ . They then consider only the  $k \ll p, n$  largest eigenvalues, and afterwards project each  $p$ -ary patient row vector  $M(j, :)$  onto  $k$  real values  $M(j, :)^T \cdot [\beta^{(1)}, \beta^{(2)}, \dots, \beta^{(k)}]$ , whose  $i^{th}$  component  $M(j, :)^T \cdot \beta^{(i)}$  depends on all component of

<sup>1</sup> Their model also includes an additive global offset  $\mu_0$ .

both the  $M(j, :)$  and  $\beta^{(i)}$  vectors. They could then use these smaller vectors, perhaps as input to a classifier.

Our ROBIC system is significantly different as it maps each patient row vector  $M(j, :)$  onto  $k$  bits (see Table 1) where the  $i^{th}$  is 1 if this patient belongs to the  $i^{th}$  bicluster, which (we will see) is true if the  $j^{th}$  component of the  $i^{th}$  eigenvector  $\beta^{(i)}$  is significantly larger than the other values in this eigenvector. N.b., this does not explicitly involve taking a dot-product, nor require using any other information from  $M(j, :)$ . (Another critical difference is that ROBIC computes these eigenvalue/eigenvectors triplets *sequentially*, while standard SVD systems compute *all* of these triplets at once; see Section 4.2.)

### 2.3 Classification Algorithms

Many researchers have developed supervised learning methods to produce classifiers that use a sample's gene expression profile to predict if that sample belongs to a known class or not. van't Veer *et al.* (2002) uses one such technique to try to predict the clinical outcome of breast cancer — *i.e.*, predict whether a patient will remain disease-free for a certain period of time. They first hand-selected 231 genes (of the 25,000 genes on the microarray) that were significantly associated with the disease outcome, *based on all samples*, then ranked these genes based on the magnitude of the correlation coefficient. They then sequentially added the top 5 remaining genes (based on this ranking), each time evaluating the quality of the current gene set by the cross-validation classification accuracy of the classifier based on that subset. They obtained an accuracy of 83% using 70 of these genes.

Unfortunately, the authors selected the features (genes) by looking at the *entire dataset*. Instead, they should have done the feature selection within each cross-validation fold. When this correction, their prediction accuracy went down by 10%, to 73% (Molla *et al.*, 2004).

We will use this Breast Cancer dataset in our studies (see  $P_1$  in Table 4 below), and also use “in-fold feature selection”.

Golub *et al.* (1999) developed another method for class prediction: They first find the genes whose expression values were strongly correlated with the class distinction to be predicted, then developed a method, “neighborhood analysis”, to define an “idealized expression pattern”, corresponding to a gene that is uniformly high in one class and uniformly low in the other. Then they use a fixed subset of “informative genes”, chosen based on their correlation with the class distinction. Each informative gene casts a weighted vote, which are summed to determine the winning class as well as a “prediction strength”. They applied their method to a human leukemia database, producing a predictor for AML versus ALL class prediction. This predictor made strong prediction for 29 of the 34 samples, and the accuracy was 100% for those 29 samples. Notice this was “training set accuracy”; we provide generalization error for an extension of this dataset,  $D_3$ .<sup>2</sup>

They tried to use their neighborhood analysis system to predict response to chemotherapy (a prognostic task), but found no informative genes, and found no evidence of a strong multi-gene expression signature correlated with clinical outcome. We report our result on this data set ( $P_2$ ) in Section 4.

Singh *et al.* (2002) also used a nearest neighbor ( $k$ -NN) approach, with gene selection, to predict the class label of prostate cancer

samples ( $D_4$ ). They then used the same approach to predict the clinical outcome, and here obtained 90% accuracy for a 5-gene model with 2 nearest neighbors.

They then used a permutation test (Mielke and Berry, 2001) to determine whether their results were statistically significant. Here, the class labels for the data points are randomly permuted — *i.e.*, randomly re-arranged, keeping only the original probability distribution of the labels. (Hence if 57% of the labels are positive in the original dataset, then 57% of the instances in the permuted set, will be labeled positive.) They then apply their learning and classification system to this new data set, with the permuted labels. Given that there is no “signal” in the labels, we expect the prediction accuracy here to be around baseline — here 57% — unless their algorithm is finding information that is not there. This is why they compare the results of their algorithm on the original unpermuted data (say 90.1% accuracy), with results of their algorithm on the permuted data, for (here) 1000 separate permuted datasets — each of which has a different permutation of the labels. In particular, they ask how often any (5-fold cross-validated) permutation score exceeded theirs, as a measure of the chance that their algorithm achieved its score on the original data by chance alone. They found accuracies that matched or exceeded their original system, 37 times. We also use their datasets ( $D_4$  and  $P_4$  below) and employ a permutation test to determine validity of our results.

Gordon *et al.* (2002) described yet another way to use gene expression measurements to predict clinical outcomes in cancer. They select 2 genes randomly, then use the ratios of their expression levels and choose thresholds to accurately distinguish between different classes of samples. This approach works fine for predicting different types of tumors in lung cancer in mesothelioma ( $D_2$ ) but when applied on the Brain clinical outcome dataset ( $P_3$  (Pomeroy *et al.*, 2002)), their prediction accuracy was only 68%.

The “shrunken centroid method” (Tibshirani *et al.*, 2002) is widely used for classifying microarray data. It shrinks the class centroids (by a shrinkage parameter  $\delta$ , typically determined by a cross-validation process) toward the overall centroids after standardizing by the within-class standard deviation for each gene. This standardization has the effect of giving higher weight to genes whose expression is stable within samples of the same class. We compare our system with this approach below.

Pranckeviciene and Somorjai (2006) argues that one should first use feature selection to identify a good, small subset of the genes, and provides a particular method, LDA EX-FS, for this task. They compare the classification results obtained with four different classifiers, using this reduced number of genes. For microarray classification, in general, they found the best classifiers are SVM and the shrunken centroids methods, especially when run on a small selection of good features. In addition, most findings confirm that feature selection leads to better classification performances, and that it is generally advisable to perform some dimensionality reduction, or at least exploratory analysis prior to classification. We borrow their ideas, of using feature selection (although of a different method, based on biclusters), and also using SVMs.

## 3 METHODS

This section describes the structure of our ROBIC algorithm for finding the biclusters, and our BIC system for using biclusters to classify microarray samples.

<sup>2</sup> In general, the  $P_i$  and  $D_j$  here refer to datasets listed in Tables 4 and 3.

### 3.1 Finding Biclusters, ROBIC

We view each bicluster as a sub-matrix  $M_{I,J}$  whose entries are “multiplicatively correlated” — that is, we identify each sample  $j \in J$  with a real value  $\beta_j$ , and each gene  $i \in I$  with a real value  $\alpha_i$ , and then view the matrix value  $m_{ij} \in M_{I,J}$  as approximately  $\alpha_i \times \beta_j$ ; see Equation 2. Perhaps genes  $I = \{g_1, g_2, g_3\}$  are correlated, in that their expressions are in a 1 to 3 to 2 ratio. The actual expression values, however, depend on the patient: here patient  $p_2$  is 10 times as responsive as patient  $p_1$ , and patient  $p_3$  (resp.,  $p_4$ ) is 2 times (resp. 4 times) as responsive as  $p_1$ . Moreover, we assume that patients 5 through  $p$  are not involved, and similarly genes 4 through  $n$  are not involved. We would then set  $\alpha' = \langle 1, 3, 2, 0, 0, \dots, 0 \rangle$  and  $\beta' = \langle 1, 10, 2, 4, 0, 0, \dots, 0 \rangle$ . (Notice we set the values for the uninvolved genes and patients to 0.) Of course, this numbering is just for pedagogical purposes; we do not anticipate that the genes in a bicluster will happen to be the first three indexed in the matrix, nor will the bicluster happen to include all-and-only the first 3 patients in the study.

Our challenge is to find these highly correlated sub-matrices. ROBIC seeks each bicluster in two steps: It first seeks the best “over-complete” vectors  $\alpha$  and  $\beta$ , which include information for the bicluster, as well as non-zero values for the other genes and samples. It then uses these vectors to find the subset of genes / samples that belong to the bicluster — as the genes/samples with the largest component value in the eigenvector.

ROBIC actually needs to find a *set* of biclusters. So after finding the appropriate  $\alpha$  and  $\beta$  vectors, it will subtract their outer product from the original matrix to produce a diminished matrix, from which it then again seeks the best bicluster, and so forth. Note that a single gene may belong in more than 1 bicluster, as might a single patient; in fact, some (patient, gene) pairs can participate in several biclusters. Also, some genes and some patients can participate in no biclusters; see Equation 4. We give details below. The next section describes BIC, which uses these biclusters to classify new patients.

Given a large  $n \times p$  matrix  $M$  (e.g.,  $50,000 \times 80$ ), ROBIC first finds two vectors,  $\alpha$  and  $\beta$ , of sizes  $|\alpha| = n$  and  $|\beta| = p$  that minimize

$$\|M - \alpha\beta^T\| = \sum_{ij} (m_{ij} - \alpha_i\beta_j)^2. \quad (3)$$

that is, we seek the best (that is, least squared error) rank-1 approximation of  $M$ . It is well known (Golub and Loan, 1989; Stewart, 1973; Leon, 1980) that these vectors correspond to the eigenvectors of the largest eigenvalue, which can be computed using standard singular value decomposition (SVD) (Golub and Loan, 1989)

$$[\alpha, v, \beta] = \text{SVD}(M, 1)$$

(Technically the  $\beta$  used in Equation 3 is  $v \times \beta$ , where  $v \in \mathfrak{R}$  is the largest eigenvalue.) We let  $\alpha \in \mathfrak{R}^n$  represent the genes’ vector and  $\beta \in \mathfrak{R}^p$  represent the samples’ vector.

ROBIC then sorts the absolute values in  $\alpha$  and  $\beta$  in descending order producing  $\alpha^{(s)}$  and  $\beta^{(s)}$  — see Figure 3[Top] — and rearranges the rows and columns in  $M$  to match this order, producing  $M_{\text{sorted}}$ .

ROBIC then applies a “hinge function” to  $\alpha^{(s)}$  (resp.,  $\beta^{(s)}$ ) to identify the subset of genes (resp., samples) that belong to this bicluster, which are the ones whose associated eigenvector components are largest.

```

ROBIC(M, K, h(·)): Bicluster_Membership
  % M is data matrix, K is # of biclusters, h(·) is hinge function
  M(1) = M.
  for k = 1..K do
    1. Compute the largest singular value decomposition of M(k),
       [α, v, β] = SVD(M(k), 1)
       • α × v × βT = best rank-1 approximation to M(k)
       • α (resp., β) represents the genes (resp., patients)
    2. α(s) = Sort(α); β(s) = Sort(β)
    3. MSorted(k) := BiSort[M(k), α/α(s), β/β(s)]
       % MSorted(k) rows are re-ordered to match to α(s)
       % and columns, reordered to match β(s)
    4. [i, j] = h(MSorted(k), α(s), β(s))
       % i (resp., j) is “hinge” value for genes (resp., patients)

    5. α(s,i) := ⟨α1(s), α2(s), ..., αi(s), 0, 0, ..., 0⟩
       β(s,j) = ⟨β1(s), β2(s), ..., βj(s), 0, 0, ..., 0⟩
       M(k+1) = MSorted(k) - α(s,i) · v · β(s,j)T
       Bk(j) = { 1 if patient j is in bicluster k
                  0 otherwise
    Return: K bicluster_memberships R = {Bk}

```

Fig. 4. Pseudo-code for rank-1 bicluster algorithm, ROBIC

It does this by finding the “best-fitting pair of lines”: That is, for each  $i = 1, \dots, n$ , find the straight line that best fits the components  $\langle \alpha_1^{(s)}, \dots, \alpha_i^{(s)} \rangle$  — i.e., that minimized the squared error — and then let  $e_1^{(i)}$  be the residual squared error. Similarly let  $e_2^{(i)}$  be the residual squared error of the best fitting line for values  $\langle \alpha_{i+1}^{(s)}, \dots, \alpha_n^{(s)} \rangle$ . ROBIC then sets

$$i^* = \text{argmin}_i \{e_1^{(i)} + e_2^{(i)}\}$$

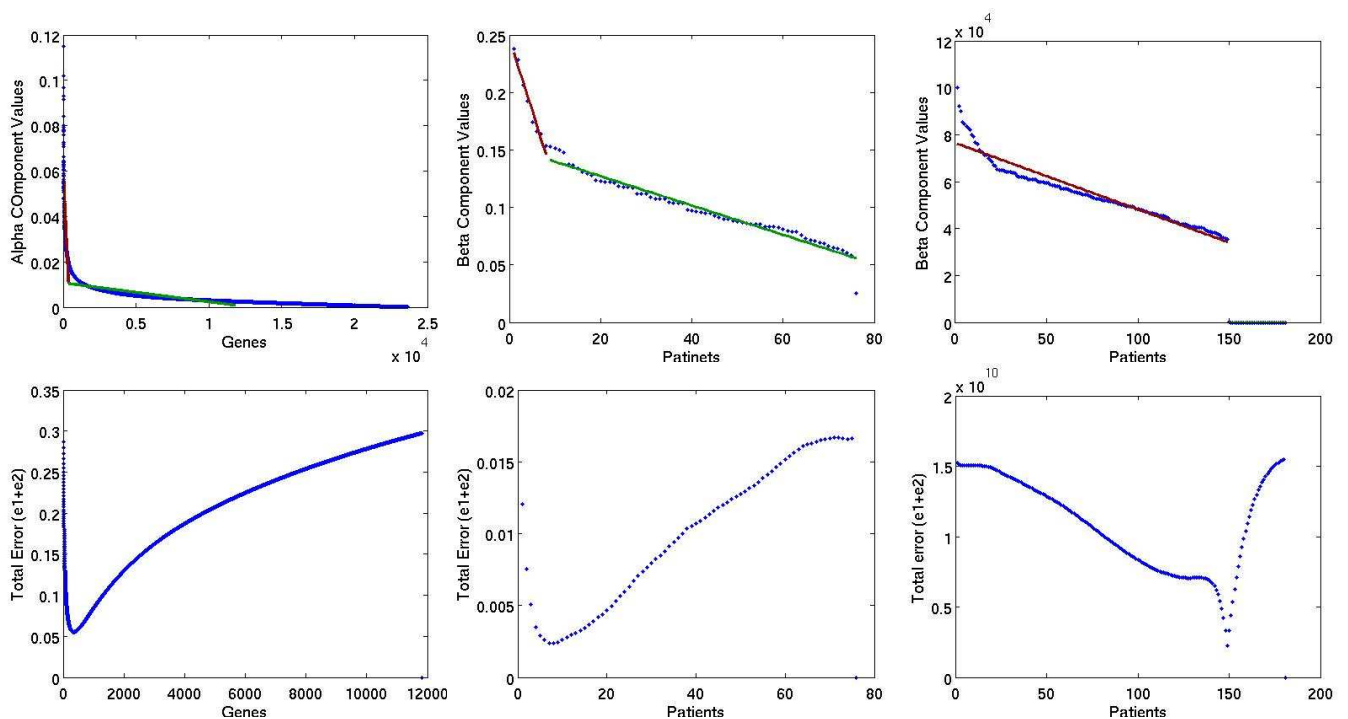
to be the index that minimizes the sum of these errors, and declares that the genes from index 1 to  $i^*$  belong to the current bicluster; see Figure 3[Bottom]. It uses a similar hinge function to determine which *samples* belong to this bicluster. (Due of the large number of genes, our implementation uses only the first half of the genes in the sorted vector  $[\alpha_1^{(s)}, \dots, \alpha_{\frac{n}{2}}^{(s)}]$  to find the best fitting lines for the genes.) It also computes a bicluster-membership vector  $B_1$ , whose  $j^{\text{th}}$  element is 1 if sample  $j$  (w.r.t. the original order) belongs to the this first cluster.

After finding the genes and samples in the first bicluster, ROBIC subtracts their values off from the data matrix, and repeats the same process on the remaining data matrix to find the next biclusters; see Equation 2.

Figure 4 summarizes our algorithm. N.b., this process does not use the class labels for the instances.

### 3.2 Bicluster Classifier, BIC

We describe each sample with an  $n$ -tuple whose  $i^{\text{th}}$  component is the gene expression value for the  $i^{\text{th}}$  gene. In the typical “learning+classification task”, we will start with  $q$  such training samples, each with a class label, and an additional  $u$  unlabeled test



**Fig. 3.** [Top] Sorted  $\alpha^{(s)}$  (resp.,  $\beta^{(s)}$ ,  $\beta^{(s)}$ ) values corresponds to genes (resp., patients, patients), with two lines superimposed. [Bottom] Error values, for each position. All figures are from the first bicluster; left 2 columns are from the Breast Cancer  $P_1$  dataset, while the far right is from Lung Cancer  $D_3$ .

instances (typically  $u = 1$ ). In general,  $p = q + u$  is much less than  $n$ . Our goal is to predict the class labels for these  $u$  test samples, based on their gene expression values.

Instead of directly using all  $n$  features, or selecting a subset of these  $n$  genes, to predict the samples' class label, we reduce the dimension of the data matrix  $M$  from  $n \times p$  real values to  $K \times p$  bits corresponding to the  $R = \{B_1, B_2, \dots, B_K\}$  matrix returned by ROBIC. Note  $K \ll n$ .

After augmenting the first  $q$  samples with their respective labels, we have data of the form shown in Table 1. BIC then takes this information, and learns a classifier based on the first  $q$  labeled samples, which it applies to the remaining  $u$  unlabeled samples, to predict their respective class labels. We experimented with several learners here, and report results using "Support Vector Machines" (SVM) and "Naive Bayes" (NB) (Alpaydin, 2004). To avoid over-fitting, it is important to use only a subset of the biclusters. We therefore use Weka's built-in in-fold feature selection algorithm to find the number of biclusters that give maximum prediction accuracy on test data; see Section 2.3. Figure 5 sketches our basic bicluster classifier algorithm, BIC. Note our system corresponds to BIC(ROBIC, ...). We will also consider other variants below.

## 4 EXPERIMENTAL RESULTS

Our goal is to produce an effective classifier of high-dimensional microarray data. To evaluate our approach, we therefore applied our BIC(ROBIC, ...) system to several publicly available datasets, and compared our results with results obtained using other approaches, including some based on an alternative type of biclusters. In each case, our ROBIC system seeks  $K = 30$  biclusters from the data set.

**Table 1.** The "bicluster membership" matrix,  $R$ , whose  $r_{jk}$  element is 1 iff the  $j^{th}$  sample is in the  $k^{th}$  bicluster, augmented with the label for the "training set" samples.

	$B_1$	$B_2$	...	$B_k$	...	$B_K$	Label Sample Class
Sample 1	1	1	...	0	...	1	+
Sample 2	1	0	...	0	...	0	-
...	...	...	...	...	...	...	...
Sample $j$	1	1	...	1	...	1	+
...	...	...	...	...	...	...	...
Sample $q$	1	1	...	0	...	1	+
Sample $q + 1$	0	0	...	0	...	1	
...	...	...	...	...	...	...	
Sample $q + u$	1	1	...	0	...	0	

### 4.1 Prognostic Datasets

We first considered datasets for predicting the clinical outcome of a certain treatment; see top "1"-labeled row of Table 4. All but one of the datasets were complete — *i.e.*, included a value for each gene-sample pair. For that one incomplete dataset ( $P_1$ , Breast Cancer), we found that 2 of the patients were missing values for roughly half (11,776) of the genes. We therefore removed those two patients.<sup>3</sup> We then removed a gene if its expression value was missing for any

<sup>3</sup> Section 4.3 discusses the ramifications of this removal.

BIC( $B(\cdot), A(\cdot), U_{Train}, U_{Test}, L_{Train}$ ):  $L_{Test}$   
 %  $B(\cdot)$  is the biclustering algorithm (e.g., ROBIC( $\cdot, 30, \dots$ ) or Plaid);  
 %  $A(\cdot)$  is the learning algorithm (e.g., SVM or NaïveBayes);  
 %  $U_{Test}$  is unlabeled test data;  $U_{Train}$  is the features associated with the training data;  
 %  $L_{Train}$  is the class labels associated with the training data; the returned  $L_{Test}$  are the predicted labels for the test data

1.  $R = B(U_{Train} \cup U_{Test}, \dots)$   
 % ... produces biclusters based on the unlabeled part of the training set and the test set. (This may involve feature selection.)
2.  $C = A(R|_{Train} + L_{Train})$   
 % ... is the classifier produced from the Training-set portion of  $R$  along with the associated training set labels
3. For  $i = 1..u$ ,  $L_{Test}[i] = C(R_{Test}[i])$   
 % Use classifier  $C$  to predict the class label of each sample in the test set.  $R_{Test}[i]$  is the  $q + i^{th}$  row of  $R$ .

Fig. 5. Pseudo-code of Bicluster Classifier algorithm BIC

Table 2. Empirical Results: Prediction accuracies of various learners, using 30 biclusters and 5-fold cross validation

	Info/Dataset	$P_1$ : Breast Cancer (van't Veer <i>et al.</i> , 2002)	$P_2$ : AML (Outcome) (Golub <i>et al.</i> , 1999)	$P_3$ : Brain (Outcome) (Pomeroy <i>et al.</i> , 2002)	$P_4$ : Prostate (Outcome) (Singh <i>et al.</i> , 2002)
1	#Samples #Genes Class distr'n	76 23,625 32 / 49 = 57.89%	15 7,129 8 / 7 = 53.33%	60 7,129 39 / 21 = 65%	21 12,600 8 / 13 = 61.90%
2	Original Result	73%		78%	<b>90%</b>
3	NaïveBayes	63.18%	46.67%	63.33%	47.62%
4	SVM	67.11%	53.33%	65%	47.62%
5	PAM	83%	60%	62%	38%
6	BIC(Plaid,NB) # of biclusters	63.16% 1	<b>73.33%</b> 19	65% 1	71.43% 6
7	BIC(Plaid,SVM) # of biclusters	63.16% 1	60% 10	65% 1	61.90% 2
8	BIC(ROBIC,NB) # of biclusters	<b>88.16%</b> 3	<b>80%</b> 30	<b>95%</b> 2	<b>76.19%</b> 1
9	BIC(ROBIC,SVM) # of biclusters	<b>90.79% ± 7.6</b> 2	<b>80% ± 18.2</b> 16	<b>95% ± 7.5</b> 2	<b>85.71% ± 12.0</b> 13
10	Permutation Test • Average • # above BIC(ROBIC,SVM) value	54.08% $\frac{0}{1000}$	51.66% $\frac{27}{1000}$	60.50% $\frac{0}{1000}$	56.77% $\frac{22}{1000}$

of the remaining patients; this meant removing 856 of the 24,481 genes.

The remaining rows of Table 4 present the prediction accuracies of various algorithms for the data sets, based on 5-fold cross-validation. The bottom of row 1 shows the “base-line” percentage, which is the accuracy of using the majority class. The second row presents the result of the original study. Rows 3 (resp., 4) show the result for applying NaïveBayes (resp., SVM) classifiers directly on the original data matrix,  $M$ . Row 5 corresponds to using shrunken centroids (PAM; (Tibshirani *et al.*, 2002)), which is one of the standard algorithms for learning classifiers for microarray data. We

empirically found setting  $\delta = 0$  worked best for the BreastCancer dataset; we therefore continued to use this setting for the other data.

Row 6 contains the results for BIC(Plaid, NB): *i.e.*, finding the biclusters using the Plaid (Lazzeroni and Owen, 2002) biclustering algorithm (rather than our ROBIC), then using NaïveBayes as the sample classification algorithm. The number under the prediction accuracy shows the number of biclusters used to obtain the specified accuracy. (Recall we selected a subset of biclusters using in-fold feature selection.) Row 7 differs from row 6 only by using SVM (not NB) for classification, BIC(Plaid,SVM). Our main results appear in rows 8 and 9, which show the prediction accuracies based on the biclusters found using our BIC(ROBIC, NB) and BIC(ROBIC, SVM)

respectively. We provide the standard deviation for the BIC(ROBIC, SVM) results. The last row in Table 4 gives the results of a 1000 trial permutation test. We see, for example, that the average accuracy of BIC(ROBIC, SVM) over 1000 permutations of the Breast Cancer dataset was 54.08%, and that none of these 1000 individual scores ever exceeded the 90.79% score on the original data. These values suggest that our AML(Outcome) results are significant with  $p \leq 0.027$  and our Prostate(Outcome), with  $p \leq 0.022$ .

## 4.2 Other Experiments

So far, focused on the difficult “prognostic” datasets. We also ran our algorithm on the typically easier “diagnostic” ones. The results appear in Table 3.<sup>4</sup> We see, again, that our BIC(ROBIC,SVM) system performs well, on these easy tasks as well — consistently obtaining 85% accuracy, or better. In many cases, the first bicluster alone is sufficient to produce these results, as it easily splits into the classes; see right-most pair of graphs in Figure 3, for Lung Cancer ( $D_3$ ).

One concern is our methodology, of finding biclusters from both the test and (the unlabeled portion of) the training data. While this is clearly appropriate in the transductive framework (Tong and Koller, 1998), it does not directly match the canonical inductive framework, which typically involves classifying a single instance at a time, based on a classifier learned from the training set.<sup>5</sup>

To match the more familiar model, we considered learning a classification function for each bicluster, based only on the training data. That is, suppose we determine that genes  $\{g_1, g_2, g_3\}$  belong to bicluster#1, which includes samples  $\{s_1, \dots, s_{30}\}$  but not  $\{s_{31}, \dots, s_{50}\}$ . We would then learn a classifier that, given the expression values of these 3 genes, predicts whether the associated sample belongs in the cluster or not.

Notice this “bicluster membership predictor”,  $BMP_1$ , was based only on the training data. Hence, the training phase would produce  $K = 30$  such  $BMP_i$ s — one for each bicluster — and also a classifier that maps from these  $K$  bits to a class label, corresponding perhaps to the disease outcome. (This is reminiscent of stacking (Wolpert, 1992).) Then, to classify a novel test sample, we would run the  $K$   $BMP_i$ s on that sample’s data to produce  $K$  bits (or fewer if used some feature selection), then let BIC use those values to predict the outcome. We did try this but found that this approach did not work well. For example, for the Breast Cancer dataset, this method produced a (5-fold cross-validated) prediction accuracy of only 73.75%.

Towards explaining this, recall that our biclusters are actually built *sequentially*: bicluster# $k$  is based on the matrix  $M^{(k)}$  formed after subtracting off biclusters numbered 1 to  $k-1$ . So a classifier learned on the values of set of genes associated with  $M^{(k)}$  will not perform well for a novel patient whose expression values correspond to  $M^{(1)}$ . Nor can we train on the  $M^{(1)}$  data, as each individual  $m_{ij}$  value is actually the sum of all biclusters that include this pair (Equation 2), which means these membership functions are not

independent — a patient’s membership in one bicluster might throw off his membership in another.

Our ROBIC system finds the top pair of eigenvectors from a matrix, then subtracts their (truncated) outer-product and recurs, finding the top pair of eigenvectors from the “reduced” matrix, and so forth,  $K = 30$  times. Another obvious approach is simply to find the top 30 pairs of eigenvectors at once — *i.e.*, simply compute  $SVD(M,30)$  once, rather than (in essence) computing  $SVD(M,1)$  30 times. We found, however, that this did not work. For the Breast Cancer dataset  $P_1$ , the (5-fold cross-validated) prediction accuracy was only 60.53%. To understand why, notice that ROBIC does not deal with the complete eigenvectors, but instead “truncates” them by setting most of the values to 0; see Step 5 of Figure 4. This non-linearity makes a significant difference.

We also experimented with a number of alternative hinge functions; see (Asgarian, 2007). None worked as well as the “best-2-line” hinge function described above.

## 4.3 Discussion

The data in Tables 4 and 3 suggest several obvious conclusions: First, we see that our BIC(ROBIC) system works extremely well: It produces at least 80% accuracy across all 8 different datasets, and over 90% in 3 of them! It is at least as good as many of the other obvious standard algorithms, and appears better than most, especially on the difficult prognostic datasets.<sup>6</sup> This is reinforced by the permutation tests, which strongly suggest that there is a signal in the bicluster membership values, which can be used in predicting prognostic and diagnostic information.

Second, we see that BIC(ROBIC, SVM) dominates BIC(ROBIC, NB), which is why we advocate using this particular classification system.<sup>7</sup>

The comparison with BIC(Plaid, ...) suggest that our ROBIC biclusters, based on “multiplicative relations between genes and samples”, can produce better classifiers than biclusters based on “additive” relations; *i.e.*, Equation 2 appears better than Equation 1 for this task.

We found that the actual classifier typically involved several different biclusters. In particular, it is not based just on the first bicluster found. This is not surprising: The claim that a set of genes within a subset of patients are correlated does not imply that this correlation is related to any pre-specified condition. So the first bicluster found might relate to some other factor (*e.g.*, gender, or some other disease or condition). Moreover, we may need to “subtract” off the influence of this dominant bicluster to reveal the bicluster that corresponds to the specific condition associated with the diagnosis or prognosis.

## 4.4 Bicluster Characteristics

Our ROBIC system found 30 biclusters for each of the 8 datasets (4 prognostic and 4 diagnostic). There were a wide range of sizes

<sup>4</sup> Greiner (2007) attempts to compare to these results to the numbers reported in the literature, and explains why this is difficult.

<sup>5</sup> In addition, our model also requires us to keep and use all of the training instances; in this way, it is like nearest neighbor (Alpaydin, 2004). Also, as the biclusters are learned independently from the labels, there is no chance that we are inadvertently “cheating” (Molla *et al.*, 2004).

<sup>6</sup> Yes, we did remove 2 of the 78 patients from the Breast Cancer dataset  $P_1$ . To be fair to the original publication, which considered all 78 patients, we could declare that our system misclassified those two patients. This would drop the accuracy of our BIC(ROBIC,SVM) system to  $90.79 \times 76/78 = 88.46$ . This is still significantly better than 73%. Notice all of the other systems listed in Table 4 were based on the set of 76 patients.

<sup>7</sup> We had earlier explored other learning algorithms in the role of SVM here, including decision trees, nearest neighbor, and SVM with various kernel, but found that SVM with the linear kernel works best.

**Table 3.** Empirical Results: Prediction accuracies for the diagnostic datasets.

Dataset Info	$D_1$ : Colon (Alon <i>et al.</i> , 1999)	$D_2$ : Lung Cancer (Gordon <i>et al.</i> , 2002)	$D_3$ : AML-ALL (Golub <i>et al.</i> , 1999)	$D_4$ : Prostate Cancer (Singh <i>et al.</i> , 2002)
#Samples	62	181	72	136
#Genes	2,000	12,533	7,129	12,600
Class distr'n	40 / 22 = 64.52%	150 / 31 = 82.87%	47 / 25 = 65.27%	77 / 59 = 56.62%
BIC(ROBIC,SVM)	88%	96.13%	84.72%	86.77%
# of biclusters	3	1	10	1

of these 240 biclusters. Figure 6 present two histograms that show the distributions of the number of patients (resp., genes) that appear in each bicluster. (To simplify the later graph, we omitted the one bicluster that included over 4500 genes.) Notice almost all of the biclusters involve a small number of patients (median 5.5), and small number of genes (median 31.5). The following table shows how many patients (resp., genes) were members of no biclusters, exactly 1 bicluster, and more than 1:

	# of patients in:	# of genes in:	
0 biclusters:	127	53563	(4)
1 bicluster:	81	18172	
>1 biclusters:	415	13010	

As expected, we see that most patients belong to many biclusters, and that the most genes participate in none. Greiner (2007) presents all of this data.

#### 4.5 Implementation and Timing

The ROBIC code (for finding the biclusters from  $M$ ) is written in Matlab. A script, written in JAVA, transforms this resulting  $R$  matrix into *Weka* (Witten and Frank, 2005) format (ARFF). We then use Weka for building the classifier on the transformed data and evaluating the results.

We ran our BIC(ROBIC, ...) system on a Pentium 4 machine with 1.3G memory. Calculating SVD( $M,1$ ) is relatively fast — only a few seconds in Matlab. The slowest part is computing the minimum error for two lines in the genes' vector, which is necessary to find the subset of genes for each bicluster. Altogether, finding each bicluster from a data matrix that has around 100 samples and 20,000 genes takes approximately 1 minute. The subsequent Weka computation required around 30 minutes, most of this time is due to the cost of the “in-fold” aspect of feature selection.

### 5 CONCLUSION

*Extensions* (1) While we focused on ways to use the biclusters to predict the class of the samples, we anticipate this approach can also be used to predict the class of the genes, *mutatis mutandis*. (2) We have only tested our system on **binary** classification tasks. We plan to extend it to more general  $r$ -ary classification tasks, as well as to regression. (3) Each bicluster is based on only a relatively small subset of the genes; it would be interesting to determine if, in fact, these genes are associated with some specific pathways. Moreover, we might be able to use these biclusters to suggest novel interactions (Madeira and Oliveira, 2004). (4) Our current system assumes complete data — that is, that we have a value for each (patient, gene) pair. We plan to explore ways to extend our system to handle incomplete data. (5) Finally, we will continue to seek other modifications

to our BIC and especially ROBIC systems, to further improve their performance.

*Contributions* DNA microarray data contain a great deal of information about a sample — often enough to predict a patient's diagnosis and prognosis. Unfortunately, as microarrays are very high-dimensional (with many 10s of thousands of genes) and very noisy, it can be very difficult to find the patterns that correspond to these predictions.

This paper presents an effective way to make accurate predictions from microarray data, by first using biclusters to reduce the dimensionality of the data. We present a novel algorithm for identifying biclusters based on the best rank-1 matrix approximation, and a method for using such biclusters to classify samples. Our empirical studies suggest that our overall system works very effectively across a number of microarray datasets; and in particular, that it works better than other standard microarray classification systems on difficult prognostic tasks. We also show that our specific “rank-1 approximation” approach is more effective than an alternative “additive” biclusterer, for this task.

See (Greiner, 2007) for more information about our implementation, and additional experimental results.

### ACKNOWLEDGMENT

We would like to thank Kathryn Graham, Ali Ghodsi, Dale Schuurmans, Mohammad Ghavamzadeh, Ray Somorjai, Volodymyr Mnih, and Csaba Szepesvari for many helpful discussions and ideas.

### REFERENCES

Alon, U., Barkai, N., Notterman, D., Gish, K., Ybarra, S., Mack, D., and Levine, A. (1999). Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. In *Proceedings of the National Academy of Sciences*, volume 96, pages 6745–6750.

Alpaydin, E. (2004). *Introduction to Machine Learning*. The MIT Press.

Alter, O., Brown, P., and Botstein, D. (2000). Singular value decomposition for genome-wide expression data processing and modeling. *PNAS*, **97**, 10101–10106.

Asgarian, N. (2007). *Rank-1 Bi-cluster Classifier*. Master's thesis, University of Alberta.

Ding, C. (2003). Unsupervised feature selection via two-way ordering in gene expression analysis. *Bioinformatics*, **19**, 1259–1266.

Golub, G. and Loan, C. V. (1989). *Matrix Computations*. The John Hopkins University Press.

Golub, T., Slonim, D., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J., Coller, H., Loh, M., Downing, J., Caligiuri, M.,



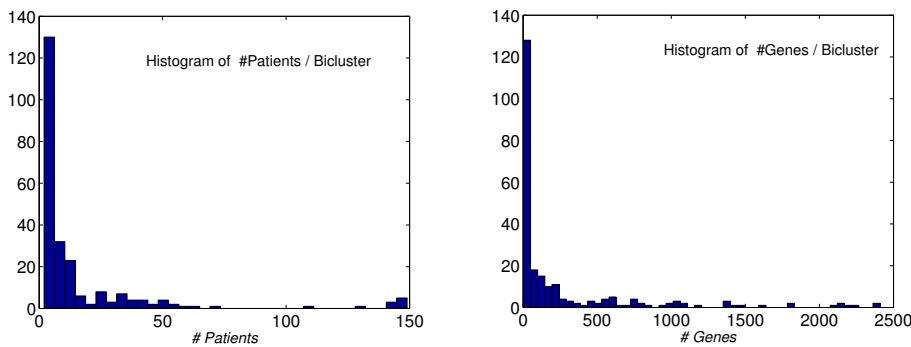


Fig. 6. Histograms showing (a) #patients/bicluster, (b) #genes/bicluster

- Bloomfield, C., and Lander, E. (1999). Molecular classification of cancer: Class discovery and class prediction by gene expression monitoring. *Science*, **286**, 531–537.
- Gordon, G. J., Jensen, R. V., Hsiao, L.-L., Gullans, S. R., Blumenstock, J. E., Ramaswamy, S., Richards, W. G., Sugarbaker, D. J., and Bueno, R. (2002). Translation of microarray data into clinically relevant cancer diagnostic tests using gene expression ratios in lung cancer and mesothelioma. *Cancer Research*, **62**, 4963–4967.
- Greiner (2007). <http://www.cs.ualberta.ca/~greiner/Research/RobiC>.
- Hastie, T., Tibshirani, R., Eisen, M., Alizadeh, A., Levy, R., Staudt, L., Chan, W., Botstein, D., and Brown, P. (2000). 'gene shaving' as a method for identifying distinct sets of genes with similar expression patterns. *Genome Biology*, **1**.
- Kluger, Y., Basri, R., Chang, J., and Gerstein, M. (2003). Spectral biclustering of microarray data: Co-clustering genes and conditions. *Genome Research*, **13**, 703–716.
- Lazzeroni, L. and Owen, A. B. (2002). Plaid models for gene expression data. *Statistica Sinica*, **12**(1), 61–86.
- Leon, S. J. (1980). *Linear Algebra With Applications*. Macmillan Publishing Co.
- Madeira, S. C. and Oliveira, A. L. (2004). Biclustering algorithms for biological data analysis: A survey. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, **1**(1), 24–45.
- McLachlan, G. J., Do, K.-A., and Ambrose, C. (2004). *Analyzing Microarray Gene Expression Data*. John Wiley and Sons, Inc., New Jersey.
- Mielke, P. W. and Berry, K. J. (2001). *Permutation Methods: A Distance Function Approach*. Springer.
- Molla, M., Waddell, M., Page, D., and Shavlik, J. (2004). Using machine learning to design and interpret gene-expression microarrays. *AI Magazine*, **25**, 23–44.
- Pomeroy, S., Tamayo, P., Gaasenbeek, M., Sturla, L., Angelo, M., McLaughlin, M., Kim, J., Goumnerova, L., Black, P., Lau, C., Allen, J., Zagzag, D., Olson, J., Curran, T., Wetmore, C., Biegel, J., Poggio, T., Mukherjee, S., Rifkin, R., Califano, A., Stolovitzky, G., Louis, D., Mesirov, J., Lander, E., and Golub, T. (2002). Prediction of central nervous system embryonal tumour outcome based on gene expression. *Letters to Nature*, **415**, 436–442.
- Pranckeviciene, E. and Somorjai, R. (2006). On classification models of gene expression microarrays: The simpler the better. In *2006 International Joint Conference on Neural Networks*, pages 6878–6885.
- Preli, A., Bleuler, S., Zimmermann, P., Wille, A., Bhlmann, P., Gruissem, W., Hennig, L., Thiele, L., and Zitzler, E. (2006). A systematic comparison and evaluation of biclustering methods for gene expression data. *Bioinformatics*, **22**(9), 1122–1129.
- Singh, D., Febbo, P., Ross, K., Jackson, D., Manola, J., Ladd, C., Tamayo, P., Renshaw, A., D'Amico, A., Richie, J., Lander, E., Loda, M., Kantoff, P., Golub, T., and Sellers, W. (2002). Gene expression correlates of clinical prostate cancer behavior. *Cancer Cell*, **1**, 203–209.
- Somorjai, R., Dolenko, B., and Baumgartner, R. (2003). Class prediction and discovery using gene microarray and proteomics mass spectroscopy data: Curses, caveats, cautions. *Bioinformatics*, **19**(12), 1484–1491.
- Stewart, G. (1973). *Introduction to Matrix Computations*. Academic Press Inc.
- Tibshirani, R., Hastie, T., Narasimhan, B., and Chu, G. (2002). Diagnosis of multiple cancer types by shrunken centroids of gene expression. In *Proceedings of the National Academy of Sciences*, volume 99, pages 6567–6572.
- Tong, S. and Koller, D. (1998). Support vector machine active learning with applications to text classification. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML-00)*, pages 287–295.
- van't Veer, L. J., Dai, H., van de Vijver, M. J., He, Y. D., Hart, A. A. M., 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). Gene expression profiling predicts clinical outcome of breast cancer. *Nature*, **415**, 530–536.
- Wall, M., Rechtsteiner, A., and Rocha, L. (2003). *Singular Value Decomposition and Principal Component Analysis*, pages 91–109. Kluwer.
- Witten, I. H. and Frank, E. (2005). *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann.
- Wolpert, D. (1992). Stacked generalization. *Neural Networks*, pages 241–260.