

Feature Selection with Conjunctions of Decision Stumps and Learning from Microarray Data

Mohak Shah, Mario Marchand, and Jacques Corbeil

Abstract

One of the objectives of designing feature selection learning algorithms is to obtain classifiers that depend on a small number of attributes *and* have verifiable future performance guarantees. There are few, if any, approaches that successfully address the two goals *simultaneously*. Performance guarantees become crucial for tasks such as microarray data analysis due to very small sample sizes resulting in limited empirical evaluation. To the best of our knowledge, such algorithms that give theoretical bounds on the future performance have not been proposed so far in the context of the classification of gene expression data. In this work, we investigate the premise of learning a conjunction (or disjunction) of *decision stumps* in Occam's Razor, Sample Compression, and PAC-Bayes learning settings for identifying a small subset of attributes that can be used to perform reliable classification tasks. We apply the proposed approaches for gene identification from DNA microarray data and compare our results to those of well known successful approaches proposed for the task. We show that our algorithm not only finds hypotheses with much smaller number of genes while giving competitive classification accuracy but also have tight risk guarantees on future performance unlike other approaches. The proposed approaches are general and extensible in terms of both designing novel algorithms and application to other domains.

Index Terms

Microarray data classification, Risk bounds, Feature selection, Gene identification.

I. INTRODUCTION

An important challenge in the problem of classification of high-dimensional data is to design a learning algorithm that can construct an accurate classifier that depends on the smallest possible number of attributes. Further, it is often desired that there be realizable guarantees associated with the future performance of such feature selection approaches. With the recent explosion in various technologies generating huge amounts of measurements, the problem of obtaining learning algorithms with performance guarantees has acquired a renewed interest.

Consider the case of biological domain where the advent of microarray technologies [Eisen and Brown, 1999, Lipshutz et al., 1999] have revolutionized the outlook on the investigation and analysis of genetic diseases. In parallel, on the classification front, many interesting results have appeared aiming to distinguish between two or more types of cells, (e.g. diseased vs. normal, or cells with different types of cancers) based on gene expression data in the case of DNA microarrays (see, for instance, [Alon et al., 1999] for results on Colon Cancer, [Golub et al., 1999] for Leukaemia). Focusing on very few genes to give insight into the class association for a microarray sample is quite important owing to a variety of reasons. For instance, a small subset of genes is easier to analyze as opposed to the set of genes output by the DNA microarray chips. It also makes it relatively easier to deduce biological relationships among them as well as study their interactions. An approach able to identify a very few number of genes can facilitate customization of chips and validation experiments— making the utilization of microarray technology cheaper, affordable, and effective.

In the view of a diseased versus a normal sample, these genes can be considered as indicators of the disease's cause. Subsequent validation study focused on these genes, their behavior, and their interactions, can lead to better understanding of the disease. Some attempts in this direction have yielded interesting results. See, for instance, a recent algorithm proposed by Wang et al. [2007] involving the identification of a gene subset based on importance ranking and subsequently combinations of genes for classification. Another example is the approach of Tibshirani et al. [2003] based on nearest shrunken centroids. Some kernel based approaches such as the BAHSIC algorithm [Song et al., 2007] and their extensions (e.g., [Shah and Corbeil, 2010] for short time-series domains) have also appeared.

The traditional methods used for classifying high-dimensional data are often characterized as either “filters” (e.g. [Furey et al., 2000, Wang et al., 2007] or “wrappers” (e.g. [Guyon et al., 2002]) depending on whether the attribute selection is performed independent of, or in conjunction with, the base learning algorithm.

M. Shah is with the Centre for Intelligent Machines, McGill University, Montreal, Canada, H3A 2A7.
E-mail: mohak@cim.mcgill.ca

M. Marchand is with the Department of Computer Science and Software Engineering, Pav. Adrien Pouliot, Laval University, Quebec, Canada, G1V-0A6.
Email: Mario.Marchand@ift.ulaval.ca

J. Corbeil is with CHUL Research Center, Laval University, Quebec (QC) Canada, G1V-4G2.
Email: Jacques.Corbeil@crchul.ulaval.ca

Despite the acceptable empirical results achieved by such approaches, there is no theoretical justification of their performance nor do they come with a guarantee on how well will they perform in the future. What is really needed is a learning algorithm that has *provably good performance guarantees* in the presence of many irrelevant attributes. This is the focus of the work presented here.

A. Contributions

The main contributions of this work come in the form of formulation of feature selection strategies within well established learning settings resulting in learning algorithms that combine the tasks of feature selection and discriminative learning. Consequently, we obtain feature selection algorithms for classification with tight realizable guarantees on their generalization error. The proposed approaches are a step towards more general learning strategies that combine feature selection with the classification algorithm *and* have tight realizable guarantees. We apply the approaches to the task of classifying microarray data where the attributes of the data sample correspond to the expression level measurements of various genes. In fact the choice of decision stumps as learning bias has in part motivated by this application. The framework is general and extensible in a variety of ways. For instance, the learning strategies proposed in this work can readily be extended to other similar tasks that can benefit from this learning bias. An immediate example would be classifying data from other microarray technologies such as in the case of Chromatin Immunoprecipitation experiments. Similarly, learning biases other than the conjunctions of decision stumps, can also be explored in the same frameworks leading to novel learning algorithms.

B. Motivation

For learning the class of conjunctions of features, we draw motivation from the guarantee that exists for this class in the following form: *if there exists a conjunction, that depends on r out of the n input attributes and that correctly classifies a training set of m examples, then the greedy covering algorithm of Haussler [1988] will find a conjunction of at most $r \ln m$ attributes that makes no training errors.* Note the absence of dependence on the number n of input attributes. The method is guaranteed to find at most $r \ln m$ attributes and, hence, depends on the number of available samples m but not on the number of attributes n to be analyzed.

We propose learning algorithms for building small conjunctions of *decision stumps*. We examine three approaches to obtain an optimal classifier based on this premise that mainly vary in the coding strategies for the threshold of each decision stump. The first two approaches attempt to do this by encoding the threshold either with message strings (Occam's Razor) or by using training examples (Sample Compression). The third strategy (PAC-Bayes) attempts to examine if an optimal classifier can be obtained by trading off the sparsity¹ of the classifier with the magnitude of the separating margin of each decision stump. In each case, we derive an upper bound on the generalization error of the classifier and subsequently use it to guide the respective algorithm. Finally, we present empirical results on the microarray data classification tasks and compare our results to the state-of-the-art approaches proposed for the task including the Support Vector Machine (SVM) coupled with feature selectors, and Adaboost. The preliminary results of this work appeared in [Marchand and Shah, 2005].

C. Organization

Section II gives the basic definitions and notions of the learning setting that we utilize and also characterizes the hypothesis class of conjunctions of decision stumps. All subsequent learning algorithms are proposed to learn this hypothesis class. Section III proposes an Occam's Razor approach to learn conjunctions of decision stumps leading to an upper bound on the generalization error in this framework. Section IV then proposes an alternate encoding strategy for the message strings using the Sample Compression framework and gives a corresponding risk bound. In Section V, we propose a PAC-Bayes approach to learn conjunction of decision stumps that enables the learning algorithm to perform an explicit non-trivial margin-sparsity trade-off to obtain more general classifiers. Section VI then proposes algorithms to learn in the three learning settings proposed in Sections III, IV and V along with a time complexity analysis. Note that the learning (optimization) strategies proposed in Section VI do not affect the respective theoretical guarantees of the learning settings. The algorithms are evaluated empirically on real world microarray datasets in Section VII. Section VIII presents a discussion on the results and also provides an analysis of the biological relevance of the selected genes in the case of each dataset, and their agreement with published findings. Finally, we conclude in Section IX.

II. DEFINITIONS

The input space \mathcal{X} consists of all n -dimensional vectors $\mathbf{x} = (x_1, \dots, x_n)$ where each real-valued component $x_i \in [A_i, B_i]$ for $i = 1, \dots, n$. Each attribute x_i for instance can refer to the expression level of gene i . Hence, A_i and B_i are, respectively, the *a priori* lower and upper bounds on values for x_i . The output space \mathcal{Y} is the set of classification labels that can be assigned to any input vector $\mathbf{x} \in \mathcal{X}$. We focus here on binary classification problems. Thus $\mathcal{Y} = \{0, 1\}$. Each example $\mathbf{z} = (\mathbf{x}, y)$ is an

¹This refers to the number of decision stumps used.

input vector $\mathbf{x} \in \mathcal{X}$ with its classification label $y \in \mathcal{Y}$ chosen *i.i.d.* from an unknown distribution D on $\mathcal{X} \times \mathcal{Y}$. The true risk $R(f)$ of any classifier f is defined as the probability that it misclassifies an example drawn according to D :

$$R(f) \stackrel{\text{def}}{=} \Pr_{(\mathbf{x}, y) \sim D} (f(\mathbf{x}) \neq y) = \mathbf{E}_{(\mathbf{x}, y) \sim D} I(f(\mathbf{x}) \neq y)$$

where $I(a) = 1$ if predicate a is true and 0 otherwise. Given a training set $S = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ of m examples, the *empirical risk* $R_S(f)$ on S , of any classifier f , is defined according to:

$$R_S(f) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{i=1}^m I(f(\mathbf{x}_i) \neq y_i) \stackrel{\text{def}}{=} \mathbf{E}_{(\mathbf{x}, y) \sim S} I(f(\mathbf{x}) \neq y)$$

The goal of any learning algorithm is to find the classifier with minimal true risk based on measuring empirical risk (and other properties) on the training sample S .

We focus on learning algorithms that construct a *conjunction of decision stumps* from a training set. Each *decision stump* is just a threshold classifier defined on a single attribute (component) x_k . More formally, a decision stump is identified by an *attribute index* $k \in \{1, \dots, n\}$, a *threshold value* $t \in \mathbb{R}$, and a *direction* $d \in \{-1, +1\}$ (that specifies whether class 1 is on the largest or smallest values of x_k). Given any input example \mathbf{x} , the output $r_{td}^k(\mathbf{x})$ of a decision stump is defined as:

$$r_{td}^k(\mathbf{x}) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } (x_k - t)d > 0 \\ 0 & \text{if } (x_k - t)d \leq 0 \end{cases}$$

We use a vector $\mathbf{k} \stackrel{\text{def}}{=} (k_1, \dots, k_{|\mathbf{k}|})$ of attribute indices $k_j \in \{1, \dots, n\}$ such that $k_1 < k_2 < \dots < k_{|\mathbf{k}|}$ where $|\mathbf{k}|$ is the number of indices present in \mathbf{k} (and thus the number of decision stumps in the conjunction)². Furthermore, We use a vector $\mathbf{t} = (t_{k_1}, t_{k_2}, \dots, t_{k_{|\mathbf{k}|}})$ of threshold values and a vector $\mathbf{d} = (d_{k_1}, d_{k_2}, \dots, d_{k_{|\mathbf{k}|}})$ of directions where $k_j \in \{1, \dots, n\}$ for $j \in \{1, \dots, |\mathbf{k}|\}$. On any input example \mathbf{x} , the output $C_{\mathbf{td}}^{\mathbf{k}}(\mathbf{x})$ of a conjunction of decision stumps is given by:

$$C_{\mathbf{td}}^{\mathbf{k}}(\mathbf{x}) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } r_{t_j d_j}^{k_j}(\mathbf{x}) = 1 \quad \forall j \in \mathbf{k} \\ 0 & \text{if } \exists j \in \mathbf{k} : r_{t_j d_j}^{k_j}(\mathbf{x}) = 0 \end{cases}$$

Finally, any algorithm that builds a conjunction can be used to build a disjunction just by exchanging the role of the positive and negative labeled examples. In order to keep our description simple, we describe here only the case of a conjunction. However, the case of disjunction follows symmetrically.

III. AN OCCAM'S RAZOR APPROACH

Our first approach towards learning the conjunction (or disjunction) of *decision stumps* is the Occam's Razor approach. Basically, we wish to obtain a hypothesis that can be coded using the least number of bits. We first propose an Occam's Razor risk bound which will ultimately guide the learning algorithm.

In the case of zero-one loss, we can model the risk of the classifier as a binomial. Let $\text{Bin}(\kappa, m, r)$ be the the binomial tail associated with a classifier of (true) risk r . Then $\text{Bin}(\kappa, m, r)$ is the probability that this classifier makes at most κ errors on a set of m examples:

$$\text{Bin}(\kappa, m, r) \stackrel{\text{def}}{=} \sum_{i=0}^{\kappa} \binom{m}{i} r^i (1-r)^{m-i}$$

The *binomial tail inversion* $\overline{\text{Bin}}(\kappa, m, \delta)$ then gives the largest risk value that a classifier can have while still having a probability of at least δ of observing at most κ errors out of m examples [Langford, 2005, Blum and Langford, 2003]:

$$\overline{\text{Bin}}(\kappa, m, \delta) \stackrel{\text{def}}{=} \sup \{r : \text{Bin}(\kappa, m, r) \geq \delta\}$$

From this definition, it follows that $\overline{\text{Bin}}(mR_S(f), m, \delta)$ is the *smallest* upper bound, which holds with probability at least $1 - \delta$, on the true risk of any classifier f with an observed empirical risk $R_S(f)$ on a test set of m examples:

$$\forall f: \Pr_{S \sim D^m} (R(f) \leq \overline{\text{Bin}}(mR_S(f), m, \delta)) \geq 1 - \delta$$

Our starting point is the Occam's razor bound of Langford [2005] which is a tighter version of the bound proposed by Blumer et al. [1987]. It is also more general in the sense that it applies to any prior distribution P over any countable class of classifiers.

²Although it is possible to use up to two decision stumps on any attribute, we limit ourselves here to the case where each attribute can be used for only one decision stump.

Theorem 1 (Langford [2005]). *For any prior distribution P over any countable class \mathcal{F} of classifiers, for any data-generating distribution D , and for any $\delta \in (0, 1]$, we have:*

$$\Pr_{S \sim D^m} \left\{ \forall f \in \mathcal{F}: R(f) \leq \overline{\text{Bin}}(mR_S(f), m, P(f)\delta) \right\} \geq 1 - \delta$$

The proof (available in [Langford, 2005]) directly follows from a straightforward union bound argument and from the fact that $\sum_{f \in \mathcal{F}} P(f) = 1$. To apply this bound for conjunctions of decision stumps we thus need to choose a suitable prior P for this class. Moreover, Theorem 1 is valid when $\sum_{f \in \mathcal{F}} P(f) \leq 1$. Consequently, we will use a *subprior* P whose sum is ≤ 1 .

In our case, decision-stumps' conjunctions are specified in terms of the discrete-valued vectors \mathbf{k} and \mathbf{d} and the continuous-valued vector \mathbf{t} . We will see below that we will use a finite-precision bit string σ to specify the set of threshold values \mathbf{t} . Let us denote by $P(\mathbf{k}, \mathbf{d}, \sigma)$ the prior probability assigned to the conjunction $C_{\sigma \mathbf{d}}^{\mathbf{k}}$ described by $(\mathbf{k}, \mathbf{d}, \sigma)$. We choose a prior of the following form:

$$P(\mathbf{k}, \mathbf{d}, \sigma) = \frac{1}{\binom{n}{|\mathbf{k}|}} p(|\mathbf{k}|) \frac{1}{2^{|\mathbf{k}|}} g_{\mathbf{k}, \mathbf{d}}(\sigma)$$

where $g_{\mathbf{k}, \mathbf{d}}(\sigma)$ is the prior probability assigned to string σ given that we have chosen \mathbf{k} and \mathbf{d} . Let $\mathcal{M}(\mathbf{k}, \mathbf{d})$ be the set of all message strings that we can use given that we have chosen \mathbf{k} and \mathbf{d} . If \mathcal{I} denotes the set of all 2^n possible attribute index vectors and $\mathcal{D}_{\mathbf{k}}$ denotes the set of all $2^{|\mathbf{k}|}$ binary direction vectors \mathbf{d} of dimension $|\mathbf{k}|$, we have that $\sum_{\mathbf{k} \in \mathcal{I}} \sum_{\mathbf{d} \in \mathcal{D}_{\mathbf{k}}} \sum_{\sigma \in \mathcal{M}(\mathbf{k}, \mathbf{d})} P(\mathbf{k}, \mathbf{d}, \sigma) \leq 1$ whenever $\sum_{d=0}^n p(d) \leq 1$ and $\sum_{\sigma \in \mathcal{M}(\mathbf{k}, \mathbf{d})} g_{\mathbf{k}, \mathbf{d}}(\sigma) \leq 1 \forall \mathbf{k}, \mathbf{d}$.

The reasons motivating this choice for the prior are the following. The first two factors come from the belief that the final classifier, constructed from the group of attributes specified by \mathbf{k} , should depend only on the number $|\mathbf{k}|$ of attributes in this group. If we have complete ignorance about the number of decision stumps the final classifier is likely to have, we should choose $p(d) = 1/(n+1)$ for $d \in \{0, 1, \dots, n\}$. However, we should choose a p that decreases as we increase d if we have reasons to believe that the number of decision stumps of the final classifier will be much smaller than n . Since this is usually our case, we propose to use:

$$p(|\mathbf{k}|) = \frac{6}{\pi^2} (|\mathbf{k}| + 1)^{-2}$$

The third factor of $P(\mathbf{k}, \mathbf{d}, \sigma)$ gives equal prior probabilities for each of the two possible values of direction d_j .

To specify the distribution of strings $g_{\mathbf{k}, \mathbf{d}}(\sigma)$, consider the problem of coding a threshold value $t \in [a, b] \subset [A, B]$ where $[A, B]$ is some predefined interval in which we are permitted to choose t and where $[a, b]$ is an interval of “equally good” threshold values.³ We propose the following diadic coding scheme for the identification of a threshold value that belongs to that interval. Let l be the number of bits that we use for the code. Then, a code of l bits specifies one value among the set Λ_l of threshold values:

$$\Lambda_l \stackrel{\text{def}}{=} \left\{ \left[1 - \frac{2j-1}{2^{l+1}} \right] A + \frac{2j-1}{2^{l+1}} B \right\}_{j=1}^{2^l}$$

We denote by A_i and B_i , the respective *a priori* minimum and maximum values that the attribute i can take. These values are obtained from the definition of data. Hence, for an attribute $i \in \mathbf{k}$, given an interval $[a_i, b_i] \subset [A_i, B_i]$ of threshold values, we take the smallest number l_i of bits such that there exists a threshold value in Λ_{l_i} that falls in the interval $[a_i, b_i]$. In that way, we will need at most $\lceil \log_2((B_i - A_i)/(b_i - a_i)) \rceil$ bits to obtain a threshold value that falls in $[a_i, b_i]$.

Hence, to specify the threshold for each decision stump $i \in \mathbf{k}$, we need to specify the number l_i of bits and a l_i -bit string s_i that identifies one of the threshold values in Λ_{l_i} . The risk bound does not depend on how we actually code σ (for some receiver). It only depends on the *a priori* probabilities we assign to each possible realization of σ . We choose the following distribution:

$$g_{\mathbf{k}, \mathbf{d}}(\sigma) \stackrel{\text{def}}{=} g_{\mathbf{k}, \mathbf{d}}(l_1, s_1, \dots, l_{|\mathbf{k}|}, s_{|\mathbf{k}|}) \tag{1}$$

$$= \prod_{i \in \mathbf{k}} \zeta(l_i) \cdot 2^{-l_i} \tag{2}$$

where:

$$\zeta(a) \stackrel{\text{def}}{=} \frac{6}{\pi^2} (a+1)^{-2} \quad \forall a \in \mathbb{N} \tag{3}$$

The sum over all the possible realizations of σ gives 1 since $\sum_{i=1}^{\infty} i^{-2} = \pi^2/6$. Note that by giving equal *a priori* probability to each of the 2^{l_i} strings s_i of length l_i , we give no preference to any threshold value in Λ_{l_i} .

The distribution ζ that we have chosen for each string length l_i has the advantage of decreasing slowly so that the risk bound does not deteriorate too rapidly as l_i increases. Other choices are clearly possible. However, note that the dominant contribution comes from the 2^{-l_i} term yielding a risk bound that depends linearly in l_i .

³By a “good” threshold value, we mean a threshold value for a decision stump that would cover many negative examples and very few positive examples (see the learning algorithm).

With this choice of prior, we have the following theorem:

Theorem 2. *Given all our previous definitions and for any $\delta \in (0, 1]$, we have:*

$$\Pr_{S \sim D^m} \left(\forall \mathbf{k}, \mathbf{d}, \sigma: R(C_{\sigma \mathbf{d}}^{\mathbf{k}}) \leq \overline{\text{Bin}} \left(m R_S(C_{\sigma \mathbf{d}}^{\mathbf{k}}), m, \frac{p(|\mathbf{k}|) g_{\mathbf{k}, \mathbf{d}}(\sigma) \delta}{\binom{n}{|\mathbf{k}|} 2^{|\mathbf{k}|}} \right) \right) \geq 1 - \delta$$

Finally, we emphasize that the risk bound of Theorem 2, used in conjunction with the distribution of messages given by $g_{\mathbf{k}, \mathbf{d}}(\sigma)$, provides a guide for choosing the optimal classifier. Note that the above risk bound suggests a non-trivial trade-off between the number of attributes and the length of the message string used to encode the classifier. Indeed the risk bound may be smaller for a conjunction having a large number of attributes with small message strings (i.e., small l_i s) than for a conjunction having a small number of attributes but with large message strings.

IV. A SAMPLE COMPRESSION APPROACH

The basic idea of the Sample compression framework [Kuzmin and Warmuth, 2007] is to obtain learning algorithms with the property that the generated classifier (with respect to some training data) can often be reconstructed with a very small subset of training examples. More formally, a learning algorithm A is said to be a *sample-compression algorithm* iff there exists a *compression function* \mathcal{C} and a *reconstruction function* \mathcal{R} such that for any training sample $S = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ (where $\mathbf{z}_i \stackrel{\text{def}}{=} (\mathbf{x}_i, y_i)$), the classifier $A(S)$ returned by A is given by:

$$A(S) = \mathcal{R}(\mathcal{C}(S)) \quad \forall S \in (\mathcal{X} \times \mathcal{Y})^m$$

For a training set S , the compression function \mathcal{C} of learning algorithm A outputs a subset \mathbf{z}_i of S , called the *compression set*, and an *information message* σ , i.e., $(\mathbf{z}_i, \sigma) = \mathcal{C}(\mathbf{z}_1, \dots, \mathbf{z}_m)$. The information message σ contains the additional information needed to reconstruct the classifier from the compression set \mathbf{z}_i . Given a training sample S , we define the compression set \mathbf{z}_i by a vector of indices \mathbf{i} such that $\mathbf{i} \stackrel{\text{def}}{=} (i_1, i_2, \dots, i_{|\mathbf{i}|})$, with $i_j \in \{1, \dots, m\} \forall j$ and $i_1 < i_2 < \dots < i_{|\mathbf{i}|}$ and where $|\mathbf{i}|$ denotes the number of indices present in \mathbf{i} .

When given an arbitrary compression set \mathbf{z}_i and an arbitrary information message σ , the *reconstruction function* \mathcal{R} of a learning algorithm A must output a classifier. The information message σ is chosen from a set $\mathcal{M}(\mathbf{z}_i)$ that consists of all the distinct messages that can be attached to the compression set \mathbf{z}_i . The existence of this reconstruction function \mathcal{R} assures that the classifier returned by $A(S)$ is *always* identified by a compression set \mathbf{z}_i and an information message σ .

In sample compression settings for learning decision stumps' conjunctions, the message string consists of the attributes and directions defined above. However, the thresholds are now specified by training examples. Hence, if we have $|\mathbf{k}|$ attributes where \mathbf{k} is the set of thresholds, the compression set consists of $|\mathbf{k}|$ training examples (one per threshold).

Our starting point is the following generic Sample Compression bound [Marchand and Sokolova, 2005]:

Theorem 3. *For any sample compression learning algorithm with a reconstruction function \mathcal{R} that maps arbitrary subsets of a training set and information messages to classifiers:*

$$\Pr_{S \sim D^m} \{ \forall \mathbf{i} \in \mathcal{I}, \sigma \in \mathcal{M}(\mathbf{z}_i): R(\mathcal{R}(\sigma, \mathbf{z}_i)) \leq \epsilon(\sigma, \mathbf{z}_i, |\mathbf{j}|) \} \geq 1 - \delta$$

where

$$\begin{aligned} \epsilon(\sigma, \mathbf{z}_i, |\mathbf{j}|) = & 1 - \exp \left(\frac{-1}{m - |\mathbf{i}| - |\mathbf{j}|} \left[\ln \binom{m}{|\mathbf{i}|} + \ln \binom{m - |\mathbf{i}|}{|\mathbf{j}|} \right. \right. \\ & \left. \left. + \ln \left(\frac{1}{P_{\mathcal{M}(\mathbf{z}_i)}(\sigma)} \right) + \ln \left(\frac{1}{\zeta(|\mathbf{i}|) \zeta(|\mathbf{j}|) \delta} \right) \right] \right) \end{aligned} \quad (4)$$

and ζ is defined by Equation 3.

Now, we need to specify the distribution of messages ($P_{\mathcal{M}(\mathbf{z}_i)}(\sigma)$) for the conjunction of decision stumps. Note that in order to specify a conjunction of decision stumps, the compression set consists of one example per decision stump. For each decision stump we have one attribute and a corresponding threshold value determined by the numerical value that this attribute takes on the training example.

The learner chooses an attribute whose threshold is identified by the associated training example. The set of these training examples form the compression set. Finally, the learner chooses a direction for each attribute.

The subset of attributes that specifies the decision stumps in our compression set \mathbf{z}_i is given by the vector \mathbf{k} defined in the previous section. Moreover, since there is one decision stump corresponding to each example in the compression set, we have $|\mathbf{i}| = |\mathbf{k}|$. Now, we assign equal probability to each possible set $|\mathbf{k}|$ of attributes (and hence thresholds) that can be selected from n attributes. Moreover, we assign equal probability over the direction that each decision stump can have $(+1, -1)$. Hence, we get the following distribution of messages:

$$P_{\mathcal{M}(\mathbf{z}_i)}(\sigma) = \binom{n}{|\mathbf{k}|}^{-1} \cdot 2^{-|\mathbf{k}|} \quad \forall \sigma \quad (5)$$

Equation 5 along with the *Sample Compression Theorem* completes the bound for the conjunction of decision stumps.

V. A PAC-BAYES APPROACH

The Occam's Razor and Sample Compression, in a sense, aim at obtaining sparse classifiers with minimum number of stumps. This sparsity is enforced by selecting the classifiers with minimal encoding of the message strings and the compression set in respective cases.

We now examine if by sacrificing this sparsity in terms of a larger separating margin around the decision boundary (yielding more confidence) can lead us to classifiers with smaller generalization error. The learning algorithm is based on the PAC-Bayes approach [McAllester, 1999] that aims at providing Probably Approximately Correct (PAC) guarantees to "Bayesian" learning algorithms specified in terms of a *prior distribution* P (before the observation of the data) and a data-dependent, *posterior distribution* Q over a space of classifiers.

We formulate a learning algorithm that outputs a stochastic classifier, called the *Gibbs Classifier* G_Q defined by a data-dependent posterior Q . Our classifier will be partly stochastic in the sense that we will formulate a posterior over the threshold values utilized by the decision stumps while still retaining the deterministic nature for the selected attributes and directions for the decision stumps.

Given an input example \mathbf{x} , the Gibbs classifier first selects a classifier h according to the posterior distribution Q and then use h to assign the label $h(\mathbf{x})$ to \mathbf{x} . The risk of G_Q is defined as the expected risk of classifiers drawn according to Q :

$$R(G_Q) \stackrel{\text{def}}{=} \mathbf{E}_{h \sim Q} R(h) = \mathbf{E}_{h \sim Q} \mathbf{E}_{(\mathbf{x}, y) \sim D} I(h(\mathbf{x}) \neq y)$$

Our starting point is the PAC-Bayes theorem [McAllester, 2003, Langford, 2005, Seeger, 2002] that provides a bound on the risk of the Gibbs classifier:

Theorem 4. *Given any space \mathcal{H} of classifiers. For any data-independent prior distribution P over \mathcal{H} , we have:*

$$\Pr_{S \sim D^m} \left(\forall Q : \text{kl}(R_S(G_Q) \| R(G_Q)) \leq \frac{\text{KL}(Q \| P) + \ln \frac{m+1}{\delta}}{m} \right) \geq 1 - \delta$$

where $\text{KL}(Q \| P)$ is the Kullback-Leibler divergence between distributions⁴ Q and P :

$$\text{KL}(Q \| P) \stackrel{\text{def}}{=} \mathbf{E}_{h \sim Q} \ln \frac{Q(h)}{P(h)}$$

and where $\text{kl}(q \| p)$ is the Kullback-Leibler divergence between the Bernoulli distributions with probabilities of success q and p :

$$\text{kl}(q \| p) \stackrel{\text{def}}{=} q \ln \frac{q}{p} + (1 - q) \ln \frac{1 - q}{1 - p}$$

This bound for the risk of Gibbs classifiers can easily be turned into a bound for the risk of Bayes classifiers B_Q over the posterior Q . B_Q basically performs a majority vote (under measure Q) of binary classifiers in \mathcal{H} . When B_Q misclassifies an example \mathbf{x} , at least half of the binary classifiers (under measure Q) misclassifies \mathbf{x} . It follows that the error rate of G_Q is at least half of the error rate of B_Q . Hence $R(B_Q) \leq 2R(G_Q)$.

In our case, we have seen that decision stump conjunctions are specified in terms of a mixture of discrete parameters \mathbf{k} and \mathbf{d} and continuous parameters \mathbf{t} . If we denote by $P_{\mathbf{k}, \mathbf{d}}(\mathbf{t})$ the probability density function associated with a prior P over the class of decision stump conjunctions, we consider here priors of the form:

$$P_{\mathbf{k}, \mathbf{d}}(\mathbf{t}) = \frac{1}{\binom{n}{|\mathbf{k}|}} p(|\mathbf{k}|) \frac{1}{2^{|\mathbf{k}|}} \prod_{j \in \mathbf{k}} \frac{I(t_j \in [A_j, B_j])}{B_j - A_j}$$

As before, we have that:

$$\sum_{\mathbf{k} \in \mathcal{I}} \sum_{\mathbf{d} \in \mathcal{D}_{\mathbf{k}}} \prod_{j \in \mathbf{k}} \int_{A_j}^{B_j} dt_j P_{\mathbf{k}, \mathbf{d}}(\mathbf{t}) = 1$$

whenever $\sum_{e=0}^n p(e) = 1$.

The factors relating to the discrete components \mathbf{k} and \mathbf{d} have the same rationale as in the case of the Occam's Razor approach. However, in the case of the threshold for each decision stumps, we now consider an explicitly continuous uniform

⁴Here $Q(h)$ denotes the probability density function associated with Q , evaluated at h .

prior. As in the Occam's Razor case, we assume each attribute value x_k to be constrained, a priori, in $[A_k, B_k]$ such that A_k and B_k are obtained from the definition of the data. Hence, we have chosen a uniform prior probability density on $[A_k, B_k]$ for each t_k such that $k \in \mathbf{k}$. This explains the last factors of $P_{\mathbf{k},\mathbf{d}}(\mathbf{t})$.

Given a training set S , the learner will choose an attribute group \mathbf{k} and a direction vector \mathbf{d} deterministically. We pose the problem of choosing the threshold in a similar manner as in the case of Occam's Razor approach of Section III with the only difference that the learner identifies the interval and selects a threshold stochastically. For each attribute $x_k \in [A_k, B_k] : k \in \mathbf{k}$, a margin interval $[a_k, b_k] \subseteq [A_k, B_k]$ is chosen by the learner. A deterministic decision stump conjunction classifier is then specified by choosing the thresholds values $t_k \in [a_k, b_k]$ uniformly. It is tempting at this point to choose $t_k = (a_k + b_k)/2 \forall k \in \mathbf{k}$ (i.e., in the middle of each interval). However, the PAC-Bayes theorem offers a better guarantee for another type of deterministic classifier as we see below.

Hence, the Gibbs classifier is defined with a posterior distribution Q having all its weight on the same \mathbf{k} and \mathbf{d} as chosen by the learner but where each t_k is uniformly chosen in $[a_k, b_k]$. The KL divergence between this posterior Q and the prior P is then given by:

$$KL(Q||P) = \ln \left(\binom{n}{|\mathbf{k}|} \cdot \frac{2^{|\mathbf{k}|}}{p(|\mathbf{k}|)} \right) + \sum_{k \in \mathbf{k}} \ln \left(\frac{B_k - A_k}{b_k - a_k} \right)$$

In this limit when $[a_k, b_k] = [A_k, B_k] \forall k \in \mathbf{k}$, it can be seen that the KL divergence between the "continuous components" of Q and P vanishes. Furthermore, the KL divergence between the "discrete components" of Q and P is small for small values of $|\mathbf{k}|$ (whenever $p(|\mathbf{k}|)$ is not too small). *Hence, this KL divergence between our choices for Q and P exhibits a tradeoff between margins ($b_k - a_k$) and sparsity (small value of $|\mathbf{k}|$) for Gibbs classifiers.* Theorem 4 suggests that the G_Q with the smallest guarantee of risk $R(G_Q)$ should minimize a non trivial combination of $KL(Q||P)$ and $R_S(G_Q)$.

The posterior Q is identified by an attribute group vector \mathbf{k} , a direction vector \mathbf{d} , and intervals $[a_k, b_k] \forall k \in \mathbf{k}$. We refine the notation for our Gibbs classifier G_Q to reflect this. Hence, we use $G_{\mathbf{ab}}^{\mathbf{kd}}$ where \mathbf{a} and \mathbf{b} are the vectors formed by the unions of a_k s and b_k s respectively. We can obtain a closed-form expression for $R_S(G_{\mathbf{ab}}^{\mathbf{kd}})$ by first considering the risk $R_{(\mathbf{x},y)}(G_{\mathbf{ab}}^{\mathbf{kd}})$ on a single example (\mathbf{x}, y) since $R_S(G_{\mathbf{ab}}^{\mathbf{kd}}) = \mathbf{E}_{(\mathbf{x},y) \sim S} R_{(\mathbf{x},y)}(G_{\mathbf{ab}}^{\mathbf{kd}})$. From our definition for Q , we find that:

$$R_{(\mathbf{x},y)}(G_{\mathbf{ab}}^{\mathbf{kd}}) = (1 - 2y) \left[\prod_{k \in \mathbf{k}} \sigma_{a_k, b_k}^{d_k}(x_k) - y \right] \quad (6)$$

where:

$$\sigma_{a,b}^d(x) \stackrel{\text{def}}{=} \begin{cases} 0 & \text{if } (x < a \text{ and } d = +1) \text{ or } (b < x \text{ and } d = -1) \\ \frac{x-a}{b-a} & \text{if } a \leq x \leq b \text{ and } d = +1 \\ \frac{b-x}{b-a} & \text{if } a \leq x \leq b \text{ and } d = -1 \\ 1 & \text{if } (b < x \text{ and } d = +1) \text{ or } (x < a \text{ and } d = -1) \end{cases}$$

Note that the expression for $R_{(\mathbf{x},y)}(G_{\mathbf{td}}^{\mathbf{k}})$ is identical to the expression for $R_{(\mathbf{x},y)}(G_{\mathbf{ab}}^{\mathbf{kd}})$ except that the piece-wise linear functions $\sigma_{a_k, b_k}^{d_k}(x_k)$ are replaced by the indicator functions $I((x_k - t_k)d_k > 0)$.

The PAC-Bayes theorem provides a risk bound for the Gibbs classifier $G_{\mathbf{ab}}^{\mathbf{kd}}$. Since the Bayes classifier $B_{\mathbf{ab}}^{\mathbf{kd}}$ just performs a majority vote under the same posterior distribution as the one used by $G_{\mathbf{ab}}^{\mathbf{kd}}$, it follows that:

$$B_{\mathbf{ab}}^{\mathbf{kd}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \prod_{k \in \mathbf{k}} \sigma_{a_k, b_k}^{d_k}(x_k) > 1/2 \\ 0 & \text{if } \prod_{k \in \mathbf{k}} \sigma_{a_k, b_k}^{d_k}(x_k) \leq 1/2 \end{cases} \quad (7)$$

Note that $B_{\mathbf{ab}}^{\mathbf{kd}}$ has an *hyperbolic* decision surface. Consequently, $B_{\mathbf{ab}}^{\mathbf{kd}}$ is not representable as a conjunction of decision stumps. There is, however, no computational difficulty at obtaining the output of $B_{\mathbf{ab}}^{\mathbf{kd}}(\mathbf{x})$ for any $\mathbf{x} \in \mathcal{X}$. We now state our main theorem:

Theorem 5. *Given all our previous definitions, for any $\delta \in (0, 1]$, and for any p satisfying $\sum_{e=0}^n p(e) = 1$, we have, with probability at least $1 - \delta$ over random draws of $S \sim D^m$:*

$$\left(\forall \mathbf{k}, \mathbf{d}, \mathbf{a}, \mathbf{b}: R(G_{\mathbf{ab}}^{\mathbf{kd}}) \leq \sup \{ \epsilon: \text{kl}(R_S(G_{\mathbf{ab}}^{\mathbf{kd}})||\epsilon) \leq \psi \} \right)$$

where

$$\psi = \frac{1}{m} \left[\ln \left(\binom{n}{|\mathbf{k}|} \cdot \frac{2^{|\mathbf{k}|}}{p(|\mathbf{k}|)} \cdot \frac{m+1}{\delta} \right) + \sum_{k \in \mathbf{k}} \ln \left(\frac{B_k - A_k}{b_k - a_k} \right) \right]$$

Furthermore: $R(B_{\mathbf{ab}}^{\mathbf{kd}}) \leq 2R(G_{\mathbf{ab}}^{\mathbf{kd}}) \quad \forall \mathbf{k}, \mathbf{d}, \mathbf{a}, \mathbf{b}$.

VI. THE LEARNING ALGORITHMS

Having proposed the theoretical frameworks attempting to obtain the optimal classifiers based on various optimization criteria, we now detail the learning algorithms for these approaches. Ideally, we would like to find a conjunction of decision stumps that minimizes the respective risk bounds for each approach. Unfortunately, this cannot be done efficiently in all cases since this problem is at least as hard as the (NP-hard) minimum set cover problem as mentioned by Marchand and Shawe-Taylor [2002]. Hence, we use a set covering greedy heuristic. It consists of choosing the decision stump i with the largest utility U_i^{SC} where:

$$U_i^{SC} = |Q_i| - p|R_i| \quad (8)$$

where Q_i is the set of negative examples covered (classified as 0) by feature i , R_i is the set of positive examples misclassified by this feature, and p is a learning parameter that gives a penalty p for each misclassified positive example. Once the feature with the largest U_i is found, we remove Q_i and R_i from the training set S and then repeat (on the remaining examples) until either no more negative examples are present or that a maximum number s of features has been reached. This heuristic was also used by Marchand and Shawe-Taylor [2002] in the context of a sample compression classifier called the set covering machine. For our sample compression approach (SC), we use the above utility function U_i^{SC} .

However, for the Occam's Razor and the PAC-Bayes approaches, we need utility functions that can incorporate the optimization aspects suggested by these approaches.

A. The Occam's Razor learning algorithm

We propose the following learning strategy for Occam's Razor learning of conjunctions of decision stumps. For a fixed l_i and η , let N be the set of negative examples and P be the set of positive examples. We start with $N' = N$ and $P' = P$. Let Q_i be the subset of N' covered by decision stump i , let R_i be the subset of P' covered by decision stump i , and let l_i be the number of bits used to code the threshold of decision stump i . We choose the decision stump i that maximizes the utility U_i^{Occam} defined as:

$$U_i^{Occam} \stackrel{\text{def}}{=} \frac{|Q_i|}{N'} - p \frac{|R_i|}{P'} - \eta \cdot l_i$$

where p is the *penalty* suffered by covering (and hence, misclassifying) a positive example and η is the cost of using l_i bits for decision stump i . Once we have found a decision stump maximizing U_i , we update $N' = N' - Q_i$ and $P' = P' - R_i$ and repeat to find the next decision stump until either $N' = \emptyset$ or the maximum number v of decision stumps has been reached (early stopping the greedy). The best values for the learning parameters p , η , and v are determined by cross-validation.

B. The PAC-Bayes Learning Algorithm

Theorem 5 suggests that the learner should try to find the Bayes classifier B_{ab}^{kd} that uses a small number of attributes (*i.e.*, a small $|k|$), each with a large separating margin ($b_k - a_k$), while keeping the empirical Gibbs risk $R_S(G_{ab}^{kd})$ at a low value. As discussed earlier, we utilize the greedy set covering heuristic for learning.

In our case, however, we need to keep the Gibbs risk on S low instead of the risk of a deterministic classifier. Since the Gibbs risk is a “soft measure” that uses the piece-wise linear functions $\sigma_{a,b}^d$ instead of the “hard” indicator functions, we cannot make use of the hard utility function of Equation 8. Instead, we need a “softer” version of this utility function to take into account covering (and erring on) an example partly. That is, a negative example that falls in the linear region of a $\sigma_{a,b}^d$ is in fact partly covered and vice versa for the positive example.

Following this observation, let \mathbf{k}' be the vector of indices of the attributes that we have used so far for the construction of the classifier. Let us first define the *covering value* $\mathcal{C}(G_{ab}^{k'd})$ of $G_{ab}^{k'd}$ by the “amount” of negative examples assigned to class 0 by $G_{ab}^{k'd}$:

$$\mathcal{C}(G_{ab}^{k'd}) \stackrel{\text{def}}{=} \sum_{(\mathbf{x}, y) \in S} (1 - y) \left[1 - \prod_{j \in \mathbf{k}'} \sigma_{a_j, b_j}^{d_j}(x_j) \right]$$

We also define the *positive-side error* $\mathcal{E}(G_{ab}^{k'd})$ of $G_{ab}^{k'd}$ as the “amount” of positive examples assigned to class 0 :

$$\mathcal{E}(G_{ab}^{k'd}) \stackrel{\text{def}}{=} \sum_{(\mathbf{x}, y) \in S} y \left[1 - \prod_{j \in \mathbf{k}'} \sigma_{a_j, b_j}^{d_j}(x_j) \right]$$

We now want to add another decision stump on another attribute, call it i , to obtain a new vector \mathbf{k}'' containing this new attribute in addition to those present in \mathbf{k}' . Hence, we now introduce the *covering contribution* of decision stump i as:

$$\begin{aligned} \mathcal{C}_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}(i) &\stackrel{\text{def}}{=} \mathcal{C}(G_{\mathbf{a}'\mathbf{b}'}^{\mathbf{k}''\mathbf{d}'} - \mathcal{C}(G_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}) \\ &= \sum_{(\mathbf{x}, y) \in S} (1 - y) \left[1 - \sigma_{a_i, b_i}^{d_i}(x_i) \right] \prod_{j \in \mathbf{k}'} \sigma_{a_j, b_j}^{d_j}(x_j) \end{aligned}$$

and the *positive-side error contribution* of decision stump i as:

$$\begin{aligned} \mathcal{E}_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}(i) &\stackrel{\text{def}}{=} \mathcal{E}(G_{\mathbf{a}'\mathbf{b}'}^{\mathbf{k}''\mathbf{d}'} - \mathcal{E}(G_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}) \\ &= \sum_{(\mathbf{x}, y) \in S} y \left[1 - \sigma_{a_i, b_i}^{d_i}(x_i) \right] \prod_{j \in \mathbf{k}'} \sigma_{a_j, b_j}^{d_j}(x_j) \end{aligned}$$

Typically, the covering contribution of decision stump i should increase its “utility” and its positive-side error should decrease it. Moreover, we want to decrease the “utility” of decision stump i by an amount which would become large whenever it has a small separating margin. Our expression for $KL(Q||P)$ suggests that this amount should be proportional to $\ln((B_i - A_i)/(b_i - a_i))$. Furthermore we should compare this margin term with the *fraction* of the remaining negative examples that decision stump i has covered (instead of the absolute amount of negative examples covered). Hence the covering contribution $\mathcal{C}_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}(i)$ of decision stump i should be divided by the amount $\mathcal{N}_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}$ of negative examples that *remains* to be covered before considering decision stump i :

$$\mathcal{N}_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}} \stackrel{\text{def}}{=} \sum_{(\mathbf{x}, y) \in S} (1 - y) \prod_{j \in \mathbf{k}'} \sigma_{a_j, b_j}^{d_j}(x_j)$$

which is simply the amount of negative examples that have been assigned to class 1 by $G_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}$. If P denotes the set of positive examples, we define the *utility* $U_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}(i)$ of *adding decision stump i to $G_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}$* as:

$$U_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}(i) \stackrel{\text{def}}{=} \frac{\mathcal{C}_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}(i)}{\mathcal{N}_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}} - p \frac{\mathcal{E}_{\mathbf{ab}}^{\mathbf{k}'\mathbf{d}}(i)}{|P|} - \eta \ln \frac{B_i - A_i}{b_i - a_i}$$

where parameter p represents the *penalty* of misclassifying a positive example and η is another parameter that controls the importance of having a large margin. These learning parameters can be chosen by cross-validation. For fixed values of these parameters, the “soft greedy” algorithm simply consists of adding, to the current Gibbs classifier, a decision stump with maximum added utility until either the maximum number v of decision stumps has been reached or all the negative examples have been (totally) covered. It is understood that, during this soft greedy algorithm, we can remove an example (\mathbf{x}, y) from S whenever it is totally covered. This occurs whenever $\prod_{j \in \mathbf{k}'} \sigma_{a_j, b_j}^{d_j}(x_j) = 0$.

Hence, we use the above utility function for the PAC-Bayes learning strategy. Note that, in the case of U_i^{PB} and U_i^{Occam} , we normalize the number of covered and erred examples so as to increase their sensitivity to the respective η terms.

1) *Time Complexity Analysis:* Let us analyze the time complexity of this algorithm for fixed p and η . For each attribute, we first sort the m examples with respect to their values for the attribute under consideration. This takes $O(m \log m)$ time. Then, we examine each potential a_i value (defined by the values of that attribute on the examples). Corresponding to each a_i , we examine all the potential b_i values (all the values greater than a_i). This gives us a time complexity of $O(m^2)$. Now if k is the largest number of examples falling into $[a_i, b_i]$, calculating the covering and error contributions and then finding the best interval $[a_i, b_i]$ takes $O(km^2)$ time. Moreover, we allow $k \in O(m)$ giving us a time complexity of $O(m^3)$ for each attribute. Finally, we do this over all the attributes. Hence, the overall time complexity of the algorithm is $O(nm^3)$. Note, however, that for microarray data, we have $n \gg m$ (hence, we can consider m^3 to be a constant). Moreover once the best stump is found, we remove the examples covered by this stump from the training set and repeat the algorithm. Now, we know that greedy algorithms of this kind have the following guarantee: if there exist r decision stumps that covers all the m examples, the greedy algorithm will find at most $r \ln(m)$ decision stumps. Since we almost always have $r \in O(1)$, the running time of the whole algorithm will almost always be $\in O(nm^3 \log(m))$. The good news is, since $n \gg m$, the time complexity of our algorithm is roughly linear in n .

2) *Fixed-Margin Heuristic:* In order to show why we prefer a uniformly distributed threshold as opposed to the one fixed at the middle of the interval $[a_i, b_i]$ for each stump i , we use an alternate algorithm that we call the fixed margin heuristic. The algorithm is similar to the one described above but with an additional parameter γ . This parameter decides a fixed margin boundary around the threshold, i.e. γ decides the length of the interval $[a_i, b_i]$. The algorithm still chooses the attribute vector \mathbf{k} , the direction vector \mathbf{d} and the vectors \mathbf{a} and \mathbf{b} . However, the a_i 's and b_i 's for each stump i are chosen such that, $|b_i - a_i| = 2\gamma$. The threshold t_i is then fixed in the middle of this interval, that is $t_i = \frac{(a_i + b_i)}{2}$. Hence, for each stump i , the interval $[a_i, b_i] = [t_i - \gamma, t_i + \gamma]$. For fixed p and γ , a similar analysis as in the previous subsection yields a time complexity of $O(nm^2 \log(m))$ for this algorithm.

Data Set			SVM	SVM+gs		SVM+rfe		Adaboost	
Name	ex	Genes	Errs	Errs	S	Errs	S	Iters	Errs
Colon	62	2000	12.8±1.4	14.4±3.5	256	15.4±4.8	128	20	15.2±2.1
B_MD	34	7129	13.2±1	7.2±2.6	32	10.4±2.4	64	20	9.8±1.1
C_MD	60	7129	28.2±2.2	23.1±2.8	1024	28.2±2.2	7129	50	21.2±2.4
Leuk	72	7129	21.3±1.4	14±2.8	64	21±3.2	256	20	17.8±1.8
Lung	52	918	8.8±1.3	6.8±1.9	64	7.2±1.8	32	1	2.4±1.4
BreastER	49	7129	15.3±2.4	10.3±2.7	256	11.2±2.8	256	50	9.8±1.7

TABLE I

RESULTS OF SVM, SVM COUPLED WITH GOLUB’S FEATURE SELECTION ALGORITHM (FILTER), SVM WITH RECURSIVE FEATURE ELIMINATION (WRAPPER) AND ADABOOST ALGORITHMS ON GENE EXPRESSION DATASETS.

Data Set			Occam			SC	
Name	ex	Genes	Errs	S	bits	Errs	S
Colon	62	2000	23.6±1.2	1.8±.6	6	18.2±1.8	1.2±.6
B_MD	34	7129	17.2±1.8	1.2±.8	3	17.2±1.3	1.4±.8
C_MD	60	7129	28.6±1.8	2.6±1.1	4	29.2±1.1	1.2±.6
Leuk	72	7129	27.8±1.7	2.2±.8	6	27.3±1.7	1.4±.7
Lung	52	918	21.7±1.1	1.8±1.2	5	18±1.3	1.2±.5
BreastER	49	7129	25.4±1.2	3.2±.6	2	21.2±1.5	1.4±.5

TABLE II

RESULTS OF THE PROPOSED OCCAM’S RAZOR AND SAMPLE COMPRESSION LEARNING ALGORITHMS ON GENE EXPRESSION DATASETS.

VII. EMPIRICAL RESULTS

The proposed approaches for learning conjunctions of *decision stumps* were tested on the six real-world binary microarray datasets viz. the *colon tumor* [Alon et al., 1999], the *Leukaemia* [Golub et al., 1999], the *B_MD* and *C_MD* Medulloblastomas data [Pomeroy et al., 2002], the *Lung* [Garber et al., 2001], and the *BreastER* data [West et al., 2001].

The *colon tumor* data set [Alon et al., 1999] provides the expression levels of 40 tumor and 22 normal colon tissues measured for 6500 human genes. We use the set of 2000 genes identified to have the highest minimal intensity across the 62 tissues. The *Leuk* data set [Golub et al., 1999] provides the expression levels of 7129 human genes for 47 samples of patients with Acute Lymphoblastic Leukemia (ALL) and 25 samples of patients with Acute Myeloid Leukemia (AML). The *B_MD* and *C_MD* data sets [Pomeroy et al., 2002] are microarray samples containing the expression levels of 7129 human genes. Data set *B_MD* contains 25 classic and 9 desmoplastic medulloblastomas whereas data set *C_MD* contains 39 medulloblastomas survivors and 21 treatment failures (non-survivors). The *Lung* dataset consists of gene expression levels of 918 genes of 52 patients with 39 Adenocarcinoma and 13 Squamous Cell Cancer [Garber et al., 2001]. This data has some missing values which were replaced by zeros. Finally, the *BreastER* dataset is the Breast Tumor data of West et al. [2001] used with Estrogen Receptor status to label the various samples. The data consists of expression levels of 7129 genes of 49 patients with 25 positive Estrogen Receptor samples and 24 negative Estrogen Receptor samples.

The number of examples and the number of genes in each data are given in the “ex” and “Genes” columns respectively under the “Data Set” tab in each table. The algorithms are referred to as “Occam” (Occam’s Razor), “SC” (Sample Compression) and “PAC-Bayes” (PAC-Bayes) in Tables II to V. They utilize the respective theoretical frameworks discussed in Sections III, IV and V along with the respective learning strategies of Section VI.

We have compared our learning algorithm with a linear-kernel soft-margin SVM trained both on all the attributes (gene expressions) and on a subset of attributes chosen by the filter method of Golub et al. [1999]. The filter method consists of ranking the attributes as function of the difference between the positive-example mean and the negative-example mean and then use only the first ℓ attributes. The resulting learning algorithm, named *SVM+gs* is the one used by Furey et al. [2000] for the same task. Guyon et al. [2002] claimed obtaining better results with the recursive feature elimination method but, as pointed out by Ambroise and McLachlan [2002], their work contained a methodological flaw. We use the SVM recursive feature elimination algorithm with this bias removed and present these results as well for comparison (referred to as “SVM+rfe” in Table I). Finally, we also compare our results with the state-of-the-art Adaboost algorithm. For this, we use the implementation in the Weka data mining software [Witten and Frank, 2005].

Each algorithm was tested over 20 random permutations of the datasets, with the 5-fold cross validation (CV) method. Each of the five training sets and testing sets was the same for all algorithms. The learning parameters of all algorithms and the gene subsets (for “SVM+gs” and “SVM+rfe”) were chosen from the training sets *only*. This was done by performing a second (nested) 5-fold CV on each training set.

For the gene subset selection procedure of SVM+gs, we have considered the first $\ell = 2^i$ genes (for $i = 0, 1, \dots, 12$) ranked according to the criterion of Golub et al. [1999] and have chosen the i value that gave the smallest 5-fold CV error on the training set. The “Errs” column under each algorithm in Tables I to III refer to the average (nested) 5-fold cross-validation error of the respective algorithm with one standard deviation two-sided confidence interval. The “bits” column in Table II

Data Set			PAC-Bayes		
Name	ex	Genes	S	G-errs	B-errs
Colon	62	2000	1.53±.28	14.68±1.8	14.65±1.8
B_MD	34	7129	1.2±.25	8.89±1.65	8.6±1.4
C_MD	60	7129	3.4±1.8	23.8±1.7	22.9±1.65
Leuk	72	7129	3.2±1.4	24.4±1.5	23.6±1.6
Lung	52	918	1.2±.3	4.4±.6	4.2±.8
BreastER	49	7129	2.6±1.1	12.8±.8	12.4±.78

TABLE III
RESULTS OF THE PAC-BAYES LEARNING ALGORITHM ON GENE EXPRESSION DATASETS.

refer to the number of bits used for the Occam’ Razor approach. The “G-errs” and the “B-errs” columns in Table III refer to the average nested 5-fold CV error of the optimal Gibbs classifier and the corresponding Bayes classifier with one standard deviation two-sided interval respectively.

For Adaboost, 10, 20, 50, 100, 200, 500, 1000 and 2000 iterations for each datasets were tried and the reported results correspond to the best obtained 5-fold CV error. The size values reported here (the “S” columns for “SVM+gs” and “SVM+rfe”, and “Itr” column for “AdaBoost” in Table I) correspond to the number of attributes (genes) selected most frequently by the respective algorithms over all the permutation runs.⁵ Choosing, by cross-validation, the number of boosting iteration is somewhat inconsistent with Adaboost’s goal of minimizing the empirical exponential risk. Indeed, to comply with Adaboost’s goal, we should choose a large-enough number of boosting rounds that assures the convergence of the empirical exponential risk to its minimum value. However, as shown by Zhang and Yu [2005], Boosting is known to overfit when the number of attributes exceeds the number of examples. This happens in the case of microarray experiments frequently where the number of genes far exceeds the number of samples, and is also the case in the datasets mentioned above. Early stopping is the recommended approach in such cases and hence we have followed the method described above to obtained the best number of boosting iterations.

Further, Table IV gives the result for a single run of the deterministic algorithm using the fixed-margin heuristic described above. Table V gives the results for the PAC-Bayes bound values for the results obtained for a single run of the PAC-Bayes algorithm on the respective microarray data sets. Recall that the PAC-Bayes bound provides a uniform upper bound on the risk of the Gibbs classifier. The column labels refer to the same quantities as above although the errors reported are over a single nested 5-fold CV run. The “Ratio” column of Table V refers to the average value of $(b_k - a_k)/(B_k - A_k)$ obtained over the decision stumps used by the classifiers over 5 testing folds and the “Bound” columns of Tables IV and V refer to the average risk bound of Theorem 5 multiplied by the total number of examples in respective data sets. Note, again, that these results are on a single permutation of the datasets and are presented just to illustrate the practicality of the risk bound and the rationale of not choosing the fixed-margin heuristic over the current learning strategy.

A. A Note on the Risk Bound

Note that the risk bounds are quite effective and their relevance should not be misconstrued by observing the results in just the current scenario. One of the most limiting factor in the current analysis is the unavailability of microarray data with larger number of examples. As the number of examples increase, the risk bound of Theorem 5 gives tighter guarantees. Consider, for instance, if the datasets for the Lung and Colon Cancer had 500 examples. A classifier with the same performance over 500 examples (i.e. with the same classification accuracy and number of features as currently) would have a bound of about 12 and 30 percent error instead of current 34.6 and 54.6 percent respectively. This only illustrates how the bound can be more effective as a guarantee when used on datasets with more examples. Similarly, a dataset of 1000 examples for Breast Cancer

⁵There were no close ties with classifiers with fewer genes.

Data Set			Stumps:PAC-Bayes(fixed margin)		
Name	ex	Genes	Size	Errors	Bound
Colon	62	2000	1	14	34
B_MD	34	7129	1	7	20
C_MD	60	7129	3	28	48
Leuk	72	7129	2	21	46
Lung	52	918	2	9	29
BreastER	49	7129	3	11	31

TABLE IV
RESULTS OF THE PAC-BAYES APPROACH WITH FIXED-MARGIN HEURISTIC ON GENE EXPRESSION DATASETS.

Data Set			Stumps:PAC-Bayes				
Name	ex	Genes	Ratio	Size	G-errs	B-errs	Bound
Colon	62	2000	0.42	1	12	11	33
B_MD	34	7129	0.10	1	7	7	20
C_MD	60	7129	0.08	5	21	20	45
Leuk	72	7129	0.002	3	22	21	48
Lung	52	918	0.12	1	3	3	18
BreastER	49	7129	0.09	2	11	11	29

TABLE V
AN ILLUSTRATION OF THE PAC-BAYES RISK BOUND ON A SAMPLE RUN OF THE PAC-BAYES ALGORITHM.

with a similar performance can have a bound of about 30 percent instead of current 63 percent. Hence, the current limitation in the practical application of the bound comes from limited data availability. As the number of examples increase, the bounds provides tighter guarantees and become more significant.

VIII. ANALYSIS

The results clearly show that even though “Occam” and “SC” are able to find sparse classifiers (with very few genes), they are not able to obtain acceptable classification accuracies. One possible explanation is that these two approaches focus on the most succinct classifier with their respective criterion. The Sample compression approach tries to minimize the number of genes used but does not take into account the magnitude of the separating margin and hence compromises accuracy. On the other hand, the Occam’s Razor approach tries to find a classifier that depends on margin *only indirectly*. Approaches based on sample compression as well as minimum description length have shown encouraging results in various domains. An alternate explanation for their suboptimal performance here can be seen in terms of extremely limited sample sizes. As a result, the gain in accuracy does not offset the cost of adding additional features in the conjunction. The PAC-Bayes approach seems to alleviate these problems by performing a significant margin-sparsity tradeoff. That is, the advantage of adding a new feature is seen in terms of a combination of the gain in both margin and the empirical risk. This can be compared to the strategy used by the regularization approaches. The classification accuracy of PAC-Bayes algorithm is competitive with the best performing classifier but has an added advantage, quite importantly, of using *very few genes*.

For the PAC-Bayes approach, we expect the Bayes classifier to generally perform better than the Gibbs classifier. This is reflected to some extent in the empirical results for Colon, C_MD and Leukaemia datasets. However, there is no means to prove that this will always be the case. It should be noted that there exist several different utility functions that we can use for each of the proposed learning approaches. We have tried some of these and reported results only for the ones that were found to be the best (and discussed in the description of the corresponding learning algorithms).

A noteworthy observation with regard to Adaboost is that the gene subset identified by this algorithm almost always include the ones found by the proposed PAC-Bayes approach for decision stumps. Most notably, the *only* gene *Cyclin D1*, a well known marker for Cancer, found for the lung cancer dataset is the most discriminating factor and is commonly found by both approaches. In both cases, the size of the classifier is almost always restricted to 1. These observations not only give insights into the absolute peaks worth investigating but also experimentally validates the proposed approaches.

Finally, many of the genes identified by the *final*⁶ PAC-Bayes classifier include some prominent markers for the corresponding diseases as detailed below.

A. Biological Relevance of the Selected Features

Table VI details the genes identified by the *final* PAC-Bayes classifier learned over each dataset after the parameter selection phase. There are some prominent markers identified by the classifier. Some of the main genes identified by the PAC-Bayes approach are the ones identified by previous studies for each disease— giving confidence in the proposed approach. Some of the discovered genes in this case include *Human monocyte-derived neutrophil-activating protein (MONAP) mRNA* in the case of Colon Cancer dataset and *oestrogen receptor* in the case of Breast Cancer data, *D79205_at-Ribosomal protein L39*, *D83542_at-Cadherin-15* and *U29195_at-NPTX2 Neuronal pentraxin II* in the case of Medulloblastomas datasets B_MD and C_MD. Other genes identified have biological relevance, for instance, the identification of *Adipsin*, *LAF-4* and *HOX1C* with regard to ALL/AML by our algorithm is in agreement with that of the findings of Chow et al. [2001], Hiwatari et al. [2003] and Lawrence and Largman [1992] respectively and the studies that followed.

Further, in the case of breast cancer, Estrogen receptors (ER) have shown to interact with BRCA1 to regulate VEGF transcription and secretion in breast cancer cells [Kawai et al., 2002]. These interactions are further investigated by Ma et al. [2005]. Further studies for ER have also been done. For instance, Moggs et al. [2005] discovered 3 putative estrogen-response elements in Keratin6 (the second gene identified by the PAC-Bayes classifier in the case of BreastER data) in the context of

⁶This is the classifier learned after choosing the best parameters using nested 5-fold CV and trained on the full dataset.

Dataset	Gene(s) identified by PAC-Bayes Classifier
Colon B_MD C_MD	<ol style="list-style-type: none"> 1. Hsa627 M26383-Human monocyte-derived neutrophil-activating protein (MONAP) mRNA 1. D79205_at-Ribosomal protein L39 1. S71824_at-Neural Cell Adhesion Molecule, Phosphatidylinositol-Linked Isoform Precursor 2. D83542_at-Cadherin-15 3. U29195_at-NPTX2 Neuronal pentraxin II 4. X73358_s_at-HAES-1 mRNA 5. L36069_at-High conductance inward rectifier potassium channel alpha subunit mRNA
Leuk	<ol style="list-style-type: none"> 1. M84526_at-DF D component of complement (adipsin) 2. U34360_at-Lymphoid nuclear protein (LAF-4) mRNA 3. M16937_at-Homo box c1 protein, mRNA
Lung BreastER	<ol style="list-style-type: none"> 1. GENE221X-IMAGE_841641-cyclin D1 (PRAD1-parathyroid adenomatosis 1) Hs82932 AA487486 1. X03635_at,X03635- class C, 20 probes, 20 in all_X03635 5885 - 6402 <p>Human mRNA for oestrogen receptor</p> <ol style="list-style-type: none"> 2. L42611_f_at, L42611- class A, 20 probes, 20 in L42611 1374-1954, <p>Homo sapiens keratin 6 isoform K6e <i>KRT6E</i> mRNA, complete cds</p>

TABLE VI
GENES IDENTIFIED BY THE *Final* PAC-BAYES CLASSIFIER

E2-responsive genes identified by microarray analysis of MDA-MD-231 cells that re-express ER α . An important role played by cytokeratins in cancer development is also widely known (see for instance Gusterson et al. [2005]).

Furthermore, the importance of *MONAP* in the case of colon cancer and *Adipsin* in the case of leukaemia data has further been confirmed by various rank based algorithms as detailed by Su et al. [2003] in the implementation of “RankGene”, a program that analyzes and ranks genes for the gene expression data using eight ranking criteria including Information Gain (IG), Gini Index (GI), Max Minority (MM), Sum Minority (SM), Twoing Rule (TR), t-statistic (TT), Sum of variances (SV) and one-dimensional Support Vector Machine (1S). In the case of Colon Cancer data, *MONAP* is identified as the top ranked gene by four of the eight criteria (IG, SV, TR, GI), second by one (SM), eighth by one (MM) and in top 50 by 1S. Similarly, in the case of Leukaemia data, *Adipsin* is top ranked by 1S, fifth by SM, seventh by IG, SV, TR, GI and MM and is in top 50 by TT. These observations provides a strong validation for our approaches.

Cyclin as identified in the case of Lung Cancer dataset is a well known marker for cell division whose perturbations are considered to be one of the major factors causing cancer [Driscoll et al., 1999, Masaki et al., 2003].

Finally, the discovered genes in the case of Medulloblastomas are important with regard to the neuronal functioning (esp. S71824, U29195 and L36039) and can have relevance for nervous system related tumors.

IX. CONCLUSION

Learning from high-dimensional data such as that from DNA microarrays can be quite challenging especially when the aim is to identify only a few attributes that characterizes the differences between two classes of data. We investigated the premise of learning conjunctions of *decision stumps* and proposed three formulations based on different learning principles. We observed that the approaches that aim solely to optimize sparsity or the message code with regard to the classifier’s empirical risk limits the algorithm in terms of its generalization performance, at least in the present case of small dataset sizes. By trading-off the sparsity of the classifier with the separating margin in addition to the empirical risk, the PAC-Bayes approach seem to alleviate this problem to a significant extent. This allows the PAC-Bayes algorithm to yield competitive classification performance while at the same time utilizing significantly fewer attributes.

As opposed to the traditional feature selection methods, the proposed approaches are accompanied by a *theoretical justification of the performance*. Moreover, the proposed algorithms *embed the feature selection as a part of the learning process* itself.⁷ Furthermore, the generalization error bounds are practical and can potentially guide the model (parameter) selection. When applied to classify DNA microarray data, the genes identified by the proposed approaches are found to be biologically significant as experimentally validated by various studies, an empirical justification that the approaches can successfully perform meaningful feature selection. Consequently, this represents a significant improvement in the direction of successful integration of machine learning approaches for use in high-throughput data to provide *meaningful, theoretically justifiable, and reliable* results. Such approaches that yield a compressed view in terms of a small number of biological markers can lead to a targeted and well focussed study of the issue of interest. For instance, the approach can be utilized in identifying gene subsets from the microarray experiments that should be further validated using focused RT-PCR techniques which are otherwise both costly and impractical to perform on the full set of genes.

Finally, as mentioned previously, the approaches presented in this wor have a wider relevance, and can have significant implications in the direction of designing theoretically justified feature selection algorithms. These are one of the few approaches that combines the feature selection with the learning process *and* provide generalization guarantees over the resulting classifiers *simultaneously*. This property assumes even more significance in the wake of limited size of microarray datasets since it limits the amount of empirical evaluation that can be reliably performed otherwise. Most natural extensions of the approaches and the learning bias proposed here would be in other similar domains including other forms of microarray experiments such as

⁷Note that Huang and Chang [2007] proposed one such approach. However, they need multiple SVM learning runs. Hence, their method basically works as a wrapper.

Chromatin Immunoprecipitation promoter arrays (chIP-Chip) and from Protein arrays. Within the same learning settings, other learning biases can also be explored such as classifiers represented by features or sets of features built on subsets of attributes.

ACKNOWLEDGMENT

This work was supported by the National Science and Engineering Research Council (NSERC) of Canada [Discovery Grant No. 122405 to MM], the Canadian Institutes of Health Research [operating grant to JC, training grant to MS while at CHUL] and the Canada Research Chair in Medical Genomics to JC.

REFERENCES

- U. Alon, N. Barkai, D.A. Notterman, K. Gish, S. Ybarra, D. Mack, and A.J. Levine. Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proc. Natl. Acad. Sci. USA*, 96(12):6745–6750, 1999.
- C. Ambrose and G. J. McLachlan. Selection bias in gene extraction on the basis of microarray gene-expression data. *Proc. Natl. Acad. Sci. USA*, 99(10):6562–6566, 2002.
- Avrim Blum and John Langford. PAC-MDL bounds. In *Proceedings of 16th Annual Conference on Learning Theory, COLT 2003*, Washington, DC, August 2003, volume 2777 of *Lecture Notes in Artificial Intelligence*, pages 344–357. Springer, Berlin, 2003.
- A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth. Occam’s razor. *Information Processing Letters*, 24:377–380, 1987.
- M. L. Chow, E. J. Moler, and I. S. Mian. Identifying marker genes in transcription profiling data using a mixture of feature relevance experts. *Physiol Genomics*, 5(2):99–111, 2001.
- B. Driscoll, S. Buckley, L. Barsky, K. Weinberg, K. D. Anderson, and D. Warburton. Abrogation of cyclin D1 expression predisposes lung cancer cells to serum deprivation-induced apoptosis. *Am J Physiol*, 276(4 Pt 1):L679–687, 1999.
- M. Eisen and P. Brown. DNA arrays for analysis of gene expression. *Methods Enzymology*, 303:179–205, 1999.
- T. S. Furey, N. Cristianini, N. Duffy, D. W. Bednarski, M. Schummer, and D. Haussler. Support vector machine classification and validation of cancer tissue samples using microarray expression data. *Bioinformatics*, 16:906–914, 2000.
- M. E. Garber, O. G. Troyanskaya, K. Schluens, S. Petersen, Z. Thaesler, M. Pacyna-Gengelbach, M. van de Rijn, G. D. Rosen, C. M. Perou, R. I. Whyte, R. B. Altman, P. O. Brown, D. Botstein, and I. Petersen. Diversity of gene expression in adenocarcinoma of the lung. *Proc. Natl. Acad. Sci. USA*, 98(24):13784–13789, 2001.
- T. R. Golub, D. K. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J. P. Mesirov, H. Coller, M. L. Loh, J. R. Downing, M. A. Caligiuri, C. D. Bloomfield, and E. S. Lander. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science*, 286(5439):531–537, 1999.
- B. A. Gusterson, D. T. Ross, V. J. Heath, and T. Stein. Basal cytokeratins and their relationship to the cellular origin and functional classification of breast cancer. *Breast Cancer Research*, 7:143–148, 2005.
- Isabelle Guyon, Jason Weston, Stephen Barnhill, and Vladimir Vapnik. Gene selection for cancer classification using support vector machines. *Machine Learning*, 46:389–422, 2002.
- D. Haussler. Quantifying inductive bias: AI learning algorithms and Valiant’s learning framework. *Artificial Intelligence*, 36:177–221, 1988.
- Mitsuteru Hiwatari, Tomohiko Taki, Takeshi Taketani, Masafumi Taniwaki, Kenichi Sugita, Mayuko Okuya, Mitsuoki Eguchi, Kohmei Ida, and Yasuhide Hayashi. Fusion of an AF4-related gene, LAF4, to MLL in childhood acute lymphoblastic leukemia with t(2;11)(q11;q23). *Oncogene*, 22(18):2851–2855, 2003.
- H. Huang and F. Chang. ESVM: Evolutionary support vector machine for automatic feature selection and classification of microarray data. *Biosystems*, 90(2):516–528, 2007.
- H. Kawai, H. Li, P. Chun, S. Avraham S, and H. K. Avraham. Direct interaction between BRCA1 and the estrogen receptor regulates vascular endothelial growth factor (VEGF) transcription and secretion in breast cancer cells. *Oncogene*, 21(50):7730–7739, 2002.
- Dima Kuzmin and Manfred K. Warmuth. Unlabeled compression schemes for maximum classes. *J. Mach. Learn. Res.*, 8:2047–2081, 2007. ISSN 1533-7928.
- John Langford. Tutorial on practical prediction theory for classification. *Journal of Machine Learning Research*, 3:273–306, 2005.
- H. J. Lawrence and C. Largman. Homeobox genes in normal hematopoiesis and leukemia. *Blood*, 80(10):2445–2453, 1992.
- R. Lipshutz, S. Fodor, T. Gingeras, and D. Lockhart. High density synthetic oligonucleotide arrays. *Nature Genetics*, 21(1 Suppl):20–24, 1999.
- Y. X. Ma, Y. Tomita, S. Fan, K. Wu, Y. Tong, Z. Zhao, L. N. Song, I. D. Goldberg, and E. M. Rosen. Structural determinants of the BRCA1 : estrogen receptor interaction. *Oncogene*, 24(11):1831–1846, 2005.
- Mario Marchand and Mohak Shah. PAC-bayes learning of conjunctions and classification of gene-expression data. In Lawrence K. Saul, Yair Weiss, and Léon Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 881–888. MIT Press, Cambridge, MA, 2005.

- Mario Marchand and John Shawe-Taylor. The set covering machine. *Journal of Machine Learning Research*, 3:723–746, 2002.
- Mario Marchand and Marina Sokolova. Learning with decision lists of data-dependent features. *Journal of Machine Learning Research*, 6:427–451, 2005.
- T. Masaki, Y. Shiratori, W. Rengifo, K. Igarashi, M. Yamagata, K. Kurokohchi, N. Uchida, Y. Miyauchi, H. Yoshiji, S. Watanabe, M. Omata, and S. Kuriyama. Cyclins and cyclin-dependent kinases: Comparative study of hepatocellular carcinoma versus cirrhosis. *Hepatology*, 37(3):534–543, 2003.
- David McAllester. PAC-Bayesian stochastic model selection. *Machine Learning*, 51:5–21, 2003. A preliminary version appeared in proceedings of COLT’99.
- David McAllester. Some PAC-Bayesian theorems. *Machine Learning*, 37:355–363, 1999.
- J. G. Moggs, T. C. Murphy, F. L. Lim, D. J. Moore, R. Stuckey, K. Antrobus, I. Kimber, and G. Orphanides. Anti-proliferative effect of estrogen in breast cancer cells that re-express ERalpha is mediated by aberrant regulation of cell cycle genes. *Journal of Molecular Endocrinology*, 34:535–551, 2005.
- S. L. Pomeroy, P. Tamayo, M. Gaasenbeek, L. M. Sturla, M. Angelo, M. E. McLaughlin, J. Y. Kim, L. C. Goumnerova, P. M. Black, C. Lau, J. C. Allen, D. Zagzag, J. M. Olson, T. Curran, C. Wetmore, J. A. Biegel, T. Poggio, S. Mukherjee, R. Rifkin, A. Califano, G. Stolovitzky, D. N. Louis, J. P. Mesirov, E. S. Lander, and T. R. Golub. Prediction of central nervous system embryonal tumour outcome based on gene expression. *Nature*, 415(6870):436–442, 2002.
- Matthias Seeger. PAC-Bayesian generalization bounds for gaussian processes. *Journal of Machine Learning Research*, 3: 233–269, 2002.
- Mohak Shah and Jacques Corbeil. A general framework for analyzing data from two short time-series microarray experiments. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, to appear, 2010. doi: <http://doi.ieeecomputersociety.org/10.1109/TCBB.2009.51>.
- L. Song, J. Bedo, K. M. Borgwardt, A. Gretton, and A. Smola. Gene selection via the BAHSIC family of algorithms. *Bioinformatics*, 23(13):490–498, 2007.
- Yang Su, T.M. Murali, Vladimir Pavlovic, Michael Schaffer, and Simon Kasif. RankGene: identification of diagnostic genes based on expression data. *Bioinformatics*, 19(12):1578–1579, 2003.
- R. Tibshirani, T. Hastie, B. Narasimhan, and G. Chu. Class prediction by nearest shrunken centroids with applications to dna microarrays. *Statistical Science*, 18:104–117, 2003.
- Lipo Wang, Feng Chu, and Wei Xie. Accurate cancer classification using expressions of very few genes. *IEEE/ACM Trans. Comput. Biol. Bioinformatics*, 4(1):40–53, 2007. ISSN 1545-5963. doi: <http://dx.doi.org/10.1109/TCBB.2007.1006>.
- M. West, C. Blanchette, H. Dressman, E. Huang, S. Ishida, R. Spang, H. Zuzan, J. A. Olson Jr, J. R. Marks, and J. R. Nevins. Predicting the clinical status of human breast cancer by using gene expression profiles. *Proc. Natl. Acad. Sci. USA*, 98(20): 11462–11467, 2001.
- Ian H. Witten and Eibe Frank. *Data Mining: Practical machine learning tools and techniques*, 2nd Ed. Morgan Kaufmann, San Francisco, 2005.
- T. Zhang and B. Yu. Boosting with early stopping: Convergence and consistency. *The Annals of Statistics*, 33:1538–1579, 2005.