

# A Bayesian Network Classification Methodology for Gene Expression Data

Paul Helman,<sup>1</sup> Robert Veroff,<sup>1</sup> Susan R. Atlas,<sup>2</sup> and Cheryl Willman<sup>3</sup>

## Abstract

We present new techniques for the application of the Bayesian network learning framework to the problem of classifying gene expression data. Our techniques address the complexities of learning Bayesian nets in several ways. First, we focus on classification and demonstrate how this reduces the Bayesian net learning problem to the problem of learning subnetworks consisting of a class label node and its set of parent genes. We then consider two different approaches to identifying parent sets which are supported by current evidence; one approach employs a simple greedy algorithm to search the universe of all genes, and a second approach develops and applies a gene selection algorithm whose results are incorporated as a prior to enable an exhaustive search for parent sets over a restricted universe of genes. Two other significant contributions are the construction of classifiers from multiple, competing Bayesian network hypotheses and algorithmic methods for normalizing and binning gene expression data in the absence of prior expert knowledge. Our classifiers are first developed under a cross validation regimen against two publicly available data sets and then validated on corresponding out-of-sample test sets. The classifiers attain a classification rate in excess of 90% on each of these out-of-sample test sets.

## 1. Introduction

The advent of high-density microarray technology for gene expression profiling on the genomic scale (Schena *et al.*, 1995; Lockhart *et al.*, 1996; DeResi *et al.*, 1997; Brown and Botstein, 1999) has opened new avenues of research in data analysis and knowledge discovery. With the huge quantities of data now being generated, the opportunities, as well as the challenges, appear almost limitless.

Recent literature explores several types of analyses of gene expression data:

- gene clustering, in which subsets of genes exhibiting similar expression patterns across *cases* (e.g., patients, experimental conditions, points of a time-series) are identified (Eisen *et al.*, 1998; Tavazoie *et al.*, 1999; Getz *et al.*, 2000; Rigoutsos *et al.*, 2000; Ben-Dor *et al.*, 2001);
- case clustering, in which sets of cases that exhibit similar gene expression patterns are identified (Alizadeh *et al.*, 2000; Getz *et al.*, 2000; Rigoutsos *et al.*, 2000; Bhattacharjee *et al.*, 2001);
- case classification, in which the value of one or more attributes external to expression data (e.g., disease subtype, treatment response, prognosis) is predicted from gene expression levels (Alon *et al.*, 1999; Golub *et al.*, 1999; Ben-Dor *et al.*, 2000; Ben-Dor *et al.*, 2001; Khan *et al.*, 2001; Ibrahim *et al.*, 2002; Pomeroy *et al.*, 2002; van't Veer *et al.*, 2002); and
- gene network reconstruction, in which models of the gene regulatory system are built (Friedman *et al.*, 1999; Murphy and Mian, 1999; Tobin *et al.*, 1999; Friedman *et al.*, 2000; D'haeseleer, 2000; Woolf and Wang, 2000; Pe'er *et al.*, 2001). This objective can be viewed as subsuming the others, provided that the external classification variables are included as nodes in the network.

---

<sup>1</sup>Computer Science Department, University of New Mexico, Albuquerque, NM 87131.

<sup>2</sup>Department of Physics and Astronomy and Center for Advanced Studies, University of New Mexico, Albuquerque, NM 87131.

<sup>3</sup>Department of Pathology and UNM Cancer Research and Treatment Center, UNM School of Medicine, University of New Mexico, Albuquerque, NM 87131.

Two factors influence a researcher's focus: the questions of interest in a given setting and the nature of the data sets available. Each of the goals sketched above is of great import, and, in fact, advances in one area often contribute to advances in the others. For example, the identification of strong gene clusters, in addition to indicating potentially significant biological relationships (e.g., co-regulation), in some instances may allow a set of genes to be collapsed into a single abstract unit, thereby reducing problem dimensionality and allowing other objectives to be more successfully addressed.

The data sets available may or may not include information to support classification. Training data that is labeled—associating with each training case the class to which it belongs—supports statistical methods for constructing a classifier. After training on a collection of labeled data, a classifier is constructed which, when presented with new query cases, predicts a class label from gene expression levels and other possibly relevant information which may be associated with a case. Without class-labeled data, genes and cases can be clustered but not classified. Often, however, an effort is made after the fact to construe biological significance for the clusters formed; the success of such clustering methods depends critically on there being a relationship between the measure of similarity used to perform clustering and actual biological similarity. Techniques that attempt to classify after training on labeled data are referred to as *supervised*, while those that do not utilize labels in training (e.g., many techniques for gene and case clustering) are known as *unsupervised*.

Additionally, various amounts of prior information (e.g., expert knowledge, such as previously known or suspected functional relationships) can be associated with gene expression data in an attempt to guide the analysis methods toward better results. Again, the amount of information available—and the degree of belief in this information—determines what information can be utilized and how it can be utilized. Little is understood regarding how such information can best be represented and applied within a rigorous and consistent framework. Such a framework will become of ever increasing importance as our biological knowledge base grows and as our objectives increase in their scope and complexity.

Our group at the University of New Mexico (UNM) is fortunate to have unusually large microarray data sets, with a substantial amount of associated clinical information. This clinical information can be utilized both as additional input and to establish classification criteria. For example, clinical history might be available that allows us to search for correlations between environmental factors and gene expression levels and, ultimately, biological manifestation (e.g., disease). In the realm of classification, we expect to have several interesting class labels to associate with our gene expression data, thus allowing us to explore a variety of supervised classification problems. Information that will be available to us includes disease absence or presence, disease type (e.g., leukemia subtypes), response to treatment, relapse / nonrelapse information, and karyotype.

Consequently, we are motivated to concentrate on the development of methodologies that can exploit the unusually rich amount of information to be associated with our gene expression data and to develop techniques particularly well suited to classification in this context. At the same time, we anticipate soon extending our objectives to include the construction of gene regulatory networks and wish also to be able to integrate in a rigorous way external information, such as prior identification of key controlling genes, causal relationships between genes, and known or hypothesized gene clusters. As is argued in the sections to follow, we believe that the mathematically grounded framework of *Bayesian networks* (Bayesian nets)—for example, Pearl (1988) and Heckerman *et al.* (1995)—satisfies these goals quite well. Furthermore, the ability of Bayesian nets to integrate prior knowledge with observational evidence potentially provides researchers with the ability to build incrementally solutions to problems of increasing scope and complexity.

The remainder of this paper is organized as follows. Section 2 briefly reviews the application of Bayesian nets in gene expression analysis, comparing our approach and objectives to some of the most successful related approaches appearing in the recent literature. Section 3 details the key mathematical elements of our approach to using Bayesian nets in classification. Section 4 presents alternative search methodologies which we have utilized in Bayesian net classifier construction. Section 5 describes our experimental design and presents a suite of results. Since we began developing and implementing our techniques prior to the production of microarray data at UNM,

the experimental results reported here are against two publicly available Affymetrix data sets:<sup>4</sup>

- MIT leukemia data (Golub *et al.*, 1999), for samples of two types, ALL and AML, of leukemia. This data set is available at [http://www-genome.wi.mit.edu/mpr/data\\_set\\_ALL\\_AML.html](http://www-genome.wi.mit.edu/mpr/data_set_ALL_AML.html).
- Princeton colon cancer data (Alon *et al.*, 1999), for normal and tumor tissue samples (available at <http://microarray.princeton.edu/oncology/affydata/index.html>).

At the time of this writing, the UNM data is becoming available, and a forthcoming paper will report on experimental results against those new data sets.

## 2. Bayesian Nets and Gene Expression Data: A Brief Overview

A Bayesian net (Pearl, 1988; Heckerman *et al.*, 1995) is a graph-based model for representing probabilistic relationships between random variables. The random variables, which may, for example, represent gene expression levels, are modeled as graph nodes; probabilistic relationships are captured by directed edges between the nodes and conditional probability distributions associated with the nodes. A Bayesian net asserts that each node is statistically independent of all its nondescendants, once the values of its parents (immediate ancestors) in the graph are known, i.e., a node  $n$ 's parents renders  $n$  and its nondescendants *conditionally independent*. It follows from these conditional independence assertions and the laws of probability that once a conditional distribution is associated with each node, specifying the probability that the node assumes a given value conditioned on the values assumed by the node's parents, a joint distribution for the entire set of random variables is uniquely determined. Algorithms and software packages (Lauritzen and Spiegelhalter, 1988; Jensen *et al.*, 1990; Shafer and Shenoy, 1990; Dawid, 1992; Decter, 1996; Madsen and Jensen, 1999; Cozman, 2001; Jensen, 2001) have been developed to help the analyst visualize and query Bayesian nets, making this a very convenient representational tool.

While Bayesian nets have found much use as a representational tool for modeling known probabilistic relationships, from the perspective of the gene expression analysis tasks of current interest, their primary utility lies in the fact that they also are a powerful learning paradigm. A body of work has evolved—see, for example, Buntine (1991), Buntine (1996), Dawid and Lauritzen (1993), Friedman and Goldszmidt (1996), Heckerman *et al.* (1995), Lam and Bacchus (1994), Pearl and Verma (1991), and Spiegelhalter *et al.* (1993)—in which statistical machine learning techniques utilize a combination of data (observations) and prior domain knowledge to direct a search for Bayesian nets which best explain the current state of knowledge embodied by these inputs. This makes Bayesian nets an attractive framework for gene expression analysis, since they can methodically hypothesize and test gene regulatory models—and other such relationships—using the rigorous methods of classical probability theory and statistics.

Not surprisingly then, others—for example, Friedman *et al.* (1999), Friedman *et al.* (2000), and Pe'er *et al.* (2001)—have successfully applied Bayesian nets to the domain of gene expression analysis. Approaches reported differ from those reported here both with respect to goals (e.g., the identification of gene relationships versus our classification objectives) and with respect to the heuristics employed in an attempt to tame the complexities of the problem. The three cited papers, for example, focus on reconstructing regulatory networks by identifying network relationships most strongly supported by the data and develop heuristics for construction of Bayesian nets that reveal such structure.

While construction of regulatory networks is an ultimate goal of our work as well, several more tractable subgoals can be identified. Our approach is to focus initially on certain of these subgoals which have high likelihood of translating into immediate clinical advances and later to synthesize our solutions to address the longer-term goals of full regulatory network reconstruction. These subgoals include, for example, the following.

---

<sup>4</sup>These sets were produced using the analysis algorithms of the Affymetrix Microarray Suite (MAS) Version 4.0. Future data sets will be based on the newer statistical algorithms provided by MAS Version 5.0. See <http://www.netaffx.com/products/mas.html>.

Classification, the prediction of a case's class (e.g., disease type, prognosis) from its gene expression profile and possibly other relevant information: The classification problem—which has not been tackled via Bayesian nets for the gene expression context in the work reported in (Friedman *et al.*, 1999; Friedman *et al.*, 2000; Pe'er *et al.*, 2001)—is the focus of the current paper. As we shall discuss, in Bayesian net terms, the construction of a classifier can be accomplished through the construction of alternative subnetworks induced by the class label node and candidate parent sets.

Construction of additional subnetworks to explore isolated regulatory (or other) relationships of interest: In addition to the class label subnetworks explored in classification, other direct, isolated, relationships are of interest as well. For example, the immediate relationship between designated genes' expression levels or the relationship between some external information—such as a chromosomal translocation—and the expression levels of some set of genes can be studied. Ultimately, causal links to clinical manifestations may be suggested.

Gene clusters: The Bayesian net framework can be used to investigate the strength of gene cluster relationships found by external means and to utilize in a variety of ways sets of highly correlated genes. For example, potentially large sets of correlated genes can be replaced in the network by a cluster module, thereby reducing the dimensionality of the space. Inter-cluster relationships and relationships between gene clusters and (for example) class labels can be investigated. Alternatively, analysis can be focused within a cluster in hopes of better understanding the co-regulatory effects exhibited by the cluster's members.

Utilization of expert knowledge: The previous point suggests that hypothesized relationships established outside the framework of a Bayesian net can be imported into a net as prior knowledge. This is just one example of a very important concept. A second example will be encountered later in this paper when we discuss how we employ external gene selection algorithms to guide the search for good parent sets. More generally, we ultimately wish to develop the means to import into the rigorous statistical learning paradigm of Bayesian nets prior knowledge that results from a variety of analyses, including those discussed in each of first three example subgoals. We believe this capability will provide researchers with a greatly enhanced ability to tackle problems (e.g., regulatory network reconstruction) whose complexities currently place them beyond the reach of Bayesian net learning techniques which start from only relatively uniformed priors.

That is, in contrast to the most widely reported approaches, our approach is to focus initially on goals more specific (e.g., classification) or more modest than learning an entire network or general subnetworks, and, by utilizing prior information (including the results of “external” algorithms), to obtain a more tractable search problem. These subproblems are important in their own right, and, in the future, we hope to be able to synthesize their solutions as we address the more comprehensive problems.

### 3. Focus on Classification

In the remainder of this paper, we focus on the goal of classification, assuming the availability of a labeled training set of cases, and we detail how we have utilized and tailored the Bayesian net framework in our solution.

#### 3.1. Bayesian Nets and Classification

Friedman *et al.* (1997) describes an approach to using Bayesian nets in classification as a way of improving upon the classification approach known as *naive Bayes* (Duda and Hart, 1973; Langley *et al.*, 1992). To our knowledge, their approach has not been applied in the context of classification from gene expression data. Our modeling differs significantly from that appearing in Friedman *et al.* (1997), as do our search procedures, scoring functions, and distribution synthesis methods. In particular, we tailor our methods to deal with problems caused by the very high dimensionality (i.e., large number of genes, each of which assumes a wide range of values)

inherent in gene expression data, while exploiting the belief that a relatively small number of genes is actually required to predict most classes of interest. We will contrast the Bayesian classifiers later in this section.

We view each gene as a random variable, with the class label as an additional random variable. The genes assume expression levels (which we shall bin into a small number of distinct values), and the label assumes values such as “cancer” or “no-cancer”, type of cancer, or response to treatment.  $\langle e \rangle$  denotes a vector of expression levels assumed by the set *genes* of all genes in a single case, and  $c_k$  denotes a value assumed by the class label. The classification problem can be stated as learning the posterior conditional distribution of the class label  $C$ , conditioned on the gene expression levels, that is, the collection of conditional probabilities

$$Pr\{C = c_k \mid genes = \langle e \rangle, \text{current knowledge}\},$$

one for each  $c_k$  and  $\langle e \rangle$  combination.

The *current knowledge* appearing in the conditioning event of the above probability generally includes both a training set of cases, and prior distributions over subsets of the random variables. In addition to a training set of cases and prior distributions, current knowledge may capture, for example, prior beliefs regarding biological mechanisms. From this perspective, classification is a problem of statistical density estimation. After viewing the training set—a sample of vectors of expression values with an associated class label, drawn from the same distribution as the query cases we later will be asked to classify—we apply laws of probability to update our priors and “learn” this common distribution. We then are able to estimate the probability that query  $q$ ’s class label  $q[C]$  is  $c_k$ , given that  $q$ ’s expression vector  $q[genes]$  is  $\langle e \rangle$ .

The main difficulty in this learning problem is that the huge dimensionality of  $\langle e \rangle$  implies that any realistically-sized sample will provide only extremely sparse coverage of the sample space. For example, even if continuous expression levels are partitioned into 2 or 3 discrete bins, each of the  $number\_of\_bins^{number\_of\_genes}$  combinations of (binned) expression levels of the several thousand genes which appear in the training data typically appears only once, and combinations in the query cases typically have not appeared at all in the training data. Consequently, estimation of the conditional distributions from simple joint frequencies observed in the sample is impossible.

We consider Bayesian nets in which each gene is a node, and the class label is an additional node having no children. Associated with each node  $n$  is a conditional distribution, a set of  $\theta_{n=v, par=\langle p \rangle} \equiv Pr\{n = v \mid Par(n) = \langle p \rangle\}$ , specifying a conditional probability for each value  $v$  of  $n$ , conditioned on each combination of values  $\langle p \rangle$  of the parents of  $n$ . Note that a Bayesian net is a pair  $(G, \Theta)$ , where  $G$  is a directed acyclic graph (DAG), and  $\Theta$  supplies a conditional probability  $\theta_{n=v, par=\langle p \rangle}$  for every node value, parent set-combination implied by  $G$ . Such a pair  $(G, \Theta)$  compactly encodes a unique joint distribution over the nodes of  $G$ ; this joint distribution  $Pr\{genes = \langle e \rangle, C = c_k\}$ , and any conditional distribution over the random variables represented by the nodes, can be recovered via various known graph traversal algorithms (Lauritzen and Spiegelhalter, 1988; Jensen *et al.*, 1990; Shafer and Shenoy, 1990; Dawid, 1992; Dechter, 1996; Madsen and Jensen, 1999; Cozman, 2001; Jensen, 2001).

If we had a fixed Bayesian net that encoded the true distribution from which each case is drawn, we could extract a classifier, namely the subgraph defined by the class label node  $C$  and its parent set  $Par(C)$ , along with the associated conditional distributions  $\theta_{C=c_k, par=\langle p \rangle} = Pr\{C = c_k \mid Par(C) = \langle p \rangle\}$ . Note that the conditional independence assertion associated with (leaf) node  $C$  implies that the classification of case  $q$  depends only on the expression levels of the genes in  $Par(C)$ , i.e., distribution  $Pr\{q[C] \mid q[genes]\}$  is identical to distribution  $Pr\{q[C] \mid q[Par(C)]\}$ . Note, in particular, that the classification does not depend on other aspects (other than the parent set of  $C$ ) of the graph structure of the Bayesian net.

Now, of course, we are not given the “true” Bayesian net. Rather, we are given a collection of training cases, and our task is to build a classifier. One approach is to attempt to learn one or more “plausible” Bayesian nets and proceed as described above. In light of the previous discussion, it is apparent that we can reduce our problem to the problem of identifying the parent sets of  $C$  in plausible Bayesian nets and of estimating the conditional probabilities associated with node  $C$  and these parent sets. In particular, once given a hypothesized parent set,

density estimation becomes far more tractable. Rather than being concerned with combinations of all the genes, we are concerned only with combinations of the parent set, and hence a training sample will generally provide much better coverage of this reduced space.

It is important to note how our modeling differs from that of a naive Bayesian classifier (Duda and Hart, 1973; Langley *et al.*, 1992) and from the generalization described in Friedman *et al.* (1997). A naive Bayesian classifier assumes independence of the attributes (genes), given the value of the class label. Under this assumption, the conditional probability  $Pr\{q[C] \mid q[genes]\}$  can be computed from the product  $\prod_{g_i \in genes} Pr\{q[g_i] \mid q[C]\}$  of the marginal conditional probabilities. The naive Bayesian model is equivalent to a Bayesian net in which no edges exist between the genes, and in which an edge exists between every gene and the class labels. We make neither assumption. Rather, we ignore the issue of what edges may exist between the genes, and compute  $Pr\{q[C] \mid q[genes]\}$  as  $Pr\{q[C] \mid q[Par(C)]\}$ , an equivalence that is valid regardless of what edges exist between the genes, provided only that  $Par(C)$  is a set of genes sufficient to render the class label conditionally independent of the remaining genes. Friedman *et al.* (1997) drops the independence assumption of a naive Bayesian classifier and attempts to learn edges between the attributes (genes, in our context), while maintaining an edge from the class label into each attribute. This approach yields good improvements over naive Bayesian classifiers in the experiments (application domains other than gene expression data) reported in Friedman *et al.* (1997). Our approach exploits a prior belief (supported by experimental results reported here and in other gene expression analyses) that for the gene expression application domain, only a small number of genes is necessary to render the class label (practically) conditionally independent of the remaining genes. This both makes learning parent sets  $Par(C)$  tractable, and generally allows the quantity  $Pr\{q[C] \mid q[Par(C)]\}$  to be well estimated from a training sample.

Our approach requires that we address the following issues, which are considered in the sections to follow.

- What does it mean for a Bayesian net to be plausible?
- What do we do with multiple plausible Bayesian nets?
- How do we find (the parent sets  $Par(C)$  in) plausible Bayesian nets?

### 3.2. Scoring the Nets

The derivations in this and the following section summarize and adapt to our context the work appearing in Heckerman *et al.* (1995), and we implicitly accept the set of assumptions made there.

Bayesian net structures are hypotheses. Each network structure  $G$  hypothesizes a collection of conditional independence assertions. Were hypothesis  $G$  true with probability 1, the assertions it encodes, plus the priors and observations  $D$ , would induce via the laws of probability a posterior distribution  $f(\Theta \mid G, D, prior)$  over the space of conditional distributions for  $G$ , where each  $\Theta$  in the space contains conditional distributions  $\theta_{n=v, par=<p>}$  for each node  $n$  in  $G$ . Of particular interest are expectations under this distribution of the form

$$E(\theta_{n=v, par=<p>} \mid G, D, prior) = \int f(\Theta \mid G, D, prior) \times \theta_{n=v, par=<p>} d\Theta,$$

as this is  $Pr\{n = v \mid Par(n) = <p>, G, D, prior\}$ . For classification, of course, the desired quantity is

$$\begin{aligned} & E(\theta_{C=c_k, par=<p>} \mid G, D, prior) \\ &= Pr\{C = c_k \mid Par(C) = <p>, G, D, prior\} \\ &= Pr\{C = c_k \mid <e>, G, D, prior\}, \end{aligned}$$

for any full expression vector  $<e>$  whose projection onto the parent set  $Par(C)$  is  $<p>$ .

In a learning context, we generally never obtain a single net structure  $G$  with certainty, but rather obtain a collection of plausible  $G_i$ . Therefore, it is desirable to employ a probabilistically-based scoring function, both to guide our exploration of nets, and to specify how to blend the distributions they induce. In a Bayesian framework, one scores how well a hypothesis  $G_i$  fits  $\{D, prior\}$  by computing

$$Pr\{D \mid G_i, prior\} = \int Pr\{D \mid \Theta\} \times f(\Theta \mid G_i, prior) d\Theta.$$

Then, from priors  $P(G_i)$  over network structures, we can obtain  $Pr\{G_i \mid D, prior\}$ . Such a scoring function is known as a *Bayesian metric*.

If we evaluated all possible structures  $G_i$  in this manner, the posterior distribution over joint distributions  $\Theta_j$  of the nodes in the networks is computed by

$$f(\Theta_j \mid D, prior) = \sum_{G_i} f(\Theta_j \mid G_i, D, prior) \times Pr\{G_i \mid D, prior\}.$$

The classification probabilities

$$Pr\{q[C] = c_k \mid q[genes] = \langle e \rangle, D, prior\}$$

of interest then are the expectations

$$E(\theta_{q[C]=c_k, q[genes]=\langle e \rangle} \mid D, prior) \tag{1}$$

under this distribution, and are obtained as a weighted sum of expectations, namely

$$\sum_{G_i} E(\theta_{q[C]=c_k, par=\langle p \rangle} \mid G_i, D, prior) \times Pr\{G_i \mid D, prior\}, \tag{2}$$

where each parent vector  $\langle p \rangle$  is the projection of  $\langle e \rangle$  onto the parent set  $par$  of  $C$  in each  $G_i$ . That is, the probability each  $G_i$  assigns to  $q[C]$  given  $q[genes]$  is weighted by the posterior  $Pr\{G_i \mid D, prior\}$ . In principle, if we could evaluate this sum over all  $G_i$  we would have an exact posterior—and hence classifier—given the current state of knowledge represented by our priors and the observed cases. The more peaked is the distribution  $Pr\{q[C] = c_k \mid q[genes] = \langle e \rangle, D, prior\}$  about its mode class  $c^*$ , the higher is the probability that the classification provided for query  $q$  is correct.

### 3.3. Computational Considerations

Our task can be viewed as approximating expression (2) by finding a set of nets whose respective contributions dominate (e.g., because they have relatively high posterior weights  $Pr\{G_i \mid D, prior\}$ ) the evaluation of this sum. Some empirical studies (Cooper and Herskovita, 1992; Heckerman *et al.*, 1995) indicate that, in a variety of contexts, a relatively small number of the nets considered (e.g., often 1) have weights large enough to materially influence the evaluation, since the weights drop off quickly as edges which represent necessary dependencies are omitted or edges which represent unnecessary dependencies are added. The experimental results reported in Section 5 explore the effect of varying the number of nets used in this approximation. One important conclusion we draw is that, in the context of high-dimensional gene expression data, the inclusion of more nets than is typical appears to yield better results. Our experiments indicate this to be the case both because the “polling” provided by a large number of nets is more accurate than that provided by a small number, and because a large number of nets often provides better coverage of the expression level combinations observed in the query cases (that is, the inclusion of more nets increases the chances that query  $q$ ’s expression levels projected onto some included parent sets have been observed in the training sample).

On the surface, the evaluation of even a single  $G$  seems a formidable task; both the expectations (1) and the Bayesian metric require an integration over potentially arbitrary distributions for  $\Theta$ . However, following the work of Heckerman *et al.* (1995), we assume that a prior distribution is specified in terms of a complete net and is Dirichlet. Intuitively, such a prior can be equated with an imaginary sample of joint observations of the random variables that represents the analyst's beliefs—both in terms of relative frequency counts (corresponding to prior probabilities) and absolute size (corresponding to degree of belief)—prior to observing the sample cases. This prior distribution on the nodes of a complete net induces on the nodes of any net a unique prior distribution consistent with a modest set of assumptions. Then, for any  $G$  and this induced prior distribution, plus a set of observed cases, the calculations reduce to a closed form.

In particular, the closed form for the expectation is

$$\begin{aligned}
& E(\theta_{n=v, par=<p>} \mid G, D, prior) \\
&= \int f(\Theta \mid G, D, prior) \times \theta_{n=v, par=<p>} d\Theta \\
&= (\alpha_{pv} + N_{pv}) / (\alpha_p + N_p),
\end{aligned} \tag{3}$$

where  $N_p$  is the number of cases observed in  $D$  in which  $Par(n) = <p>$ ,  $N_{pv}$  is the number of cases observed in  $D$  in which  $Par(n) = <p>$  and  $n = v$ , and  $\alpha_p$  and  $\alpha_{pv}$  are derived from prior probabilities for these combinations of values and, under our prior assignments, are extremely small (see Section 3.4 and Heckerman *et al.* (1995)). The closed form for the Bayesian metric is

$$\begin{aligned}
& Pr\{D \mid G, prior\} \\
&= \int Pr\{D \mid \Theta\} \times f(\Theta \mid G, prior) d\Theta \\
&= \prod_n \prod_p \frac{\Gamma(\alpha_p)}{\Gamma(\alpha_p + N_p)} \prod_v \frac{\Gamma(\alpha_{pv} + N_{pv})}{\Gamma(\alpha_{pv})},
\end{aligned}$$

where

$\Gamma$  is the Gamma function;

$n$  ranges over the nodes in  $G$ ;

$p$  ranges over values  $<p>$  of  $Par(n)$  for the node  $n$  fixed by the outermost  $\prod$ ;

$v$  ranges over the values of the node  $n$  fixed by the outermost  $\prod$ ; and

$\alpha_p, \alpha_{pv}, N_p, N_{pv}$  are as defined above, with respect to the node  $n$  fixed by the outermost  $\prod$ .

The above expression for  $Pr\{D \mid G, prior\}$ , which assumes a Dirichlet prior, is known as the *BD (Bayesian-Dirichlet) metric*. (Technically, the BD metric is more commonly defined in terms of the joint posterior probability  $Pr\{D, G \mid prior\}$ , which is simply the above expression multiplied by the network prior  $P(G)$ .)

Further simplifying the computational task is the observation that the scoring function is decomposable; it can be expressed as the product of scores over the nodes, where a node's score depends only on its parent set. In our restricted context of classification, this means we can ignore the score of every node except the label, effectively using the *BD* metric as an evaluator of potential parent sets. More precisely, the *BD* evaluation of a parent set  $Par(C)$  is node  $C$ 's contribution to the *BD* score of *any* Bayesian net containing this subgraph. In particular (in contrast to a naive Bayesian classifier, in which there must be no edges between genes), the decomposability of

the *BD* score allows the hypothesis represented by parent set  $Par(C)$  to be evaluated in isolation of the question of what other edges may exist in the network. Similarly, since the expectation of interest depends only on frequencies of node  $C$  and of its parent set, the remainder of the network can be ignored in our context.

### 3.4. Specification of Priors

In each of the experiments reported, we choose an uninformed prior over the distributions that can be associated with any given network structure. In particular, we employ an extremely small equivalent sample size (Heckerman *et al.*, 1995) of 0.001, and assign each joint combination of variable values equal probability. There then is a simple translation of this prior to priors over the possible conditional distributions in any given network structure, yielding the  $\alpha_{pv}$  and  $\alpha_p$  appearing in expression (3). Our choice of prior minimizes its impact on posterior calculations, allowing the data to dominate.

The network structures  $G$  are assigned a uniform prior also, but after various prunings (see Section 4) have been imposed. In the context of our minimal-knowledge greedy algorithm, a prior which assigns equal probability to each DAG in which the class label has  $\mathcal{M}$  or fewer parents (and zero probability to all other DAGs) is used, for some specified maximum cardinality choice  $\mathcal{M}$ . In the context of the external gene selection algorithms, a prior which assigns equal probability to each DAG in which the class label has  $\mathcal{M}$  or fewer parents, each of which is a member of the selected set of genes (and zero probability to all other DAGs), is used.

Current research is considering how various types of expert biological information can be incorporated into priors and utilized by our methods. This is an area we believe to be critically important to future advances.

### 3.5. Binning Issues

Though Bayesian nets can be utilized to represent continuous distributions, most Bayesian net procedures assume that the random variables take on only a small number (e.g., 2 or 3) of discrete values. This requires procedures to discretize (i.e., collapse) typically continuous gene expression values. We describe in Section 4 the two relatively simple approaches we have used with our current search procedures. The first method bins expression values into “low”, “medium”, and “high” based on the distance of a particular expression value from the gene’s mean expression value. The second method is more closely coupled with our external gene selection method and produces a binary binning based on a maximal “point of separation” in the training data between the classes.

While these simple methods have produced good classification results, we point out here that there are many interesting avenues of research in which the binning procedure is more integrated with the search for good Bayesian nets, and candidate binnings are evaluated in the same framework as are other aspects of the nets. We consider this to be an important avenue for future research.

### 3.6. A Multi-Parent-Set Classifier

We have indicated how a parent set of the class label corresponds to the relevant (for classification) subgraph of a Bayesian net and, with expression (2), how the class distributions associated with each parent set in a collection of parent sets are combined by means of the *BD* scoring metric. Our method then is to build a classifier from some number  $\mathcal{PS}$  of parent sets that score high under the *BD* metric. That is, we perform some form of search (see the next section), selecting the  $\mathcal{PS}$  top scoring parent sets, and these are the sets whose distributions contribute the terms for our approximation of the expression (2). We see from expression (3) that the individual probabilities contributed are simply of the form  $(\alpha_{pv} + N_{pv})/(\alpha_p + N_p)$ .

An important phenomenon results from the sparseness of the data, especially in the high dimensional space of microarray data. It is possible that the combinations of values appearing in  $q[par_i]$  for some of the parent sets  $par_i$  are not seen in training or seen only minimally (for example, one or two occurrences). The distributions

yielded by such nets will then reflect only the prior, which (as we shall generally assume) is uninformed, yielding equal class probabilities, or will be determined by the handful of training cases with this  $par_i$  combination. It is important to note that this is the correct posterior distribution under the hypothesis of this parent set and given current knowledge and should not be interpreted as a “weak” or “missing” distribution simply because it is based on a small or empty sample. The strength of this distribution as it contributes to (2) is determined solely by the  $BD$  fit. A dispersed distribution (e.g., uniform) learned from a small sample and a peaked distribution learned from a large sample contribute their expectation in the same way, their relative contributions to the posterior affected only by their  $BD$  fit.<sup>5</sup>

Is it correct to treat the sparse-sample based distributions on equal footing with large-sample based distributions? We consider the variance of the distribution. Variance reflects, among other characteristics, how much the distribution may be expected to change if more data is observed. In the case of high variance, it is not unlikely that new data will shift the distribution dramatically.

The variance of the posterior  $Pr\{C = c_k | Par(C) = \langle p \rangle, G, D, prior\}$  of a binary-valued class label, being a Dirichlet distribution, is

$$(Pr\{C = c_k | Par(C) = \langle p \rangle\} \times (1 - Pr\{C = c_k | Par(C) = \langle p \rangle\})) / (\alpha_p + N_p + 1).$$

So, an interpretation is, when the “sample size”  $N_p$  is small, or when the probability is spread evenly across the classes, variance is relatively high, and the distribution is possibly “unstable” in the presence of additional observations. While the posterior distribution it yields is undeniable given the current state of knowledge, it is not unlikely to change dramatically given new data. In this sense, it is less “reliable”.

We have experimented with two heuristics for adjusting a parent set’s contribution to the evaluation of a query case in order to address the issue of the variance of the distribution. Note that unlike a set’s  $BD$  score, which is used in parent set selection as well as for a weight in the posterior computation (2), this adjustment is *query specific*, reflecting the amount of variance  $var(q)$  in the distribution of a particular query  $q$ ’s (unknown) label. The two adjustments considered are:

- When evaluating a query  $q$ , set to zero the weight in (2) of any parent set  $par_i$  such that  $q[par_i]$  has no occurrences in the training sample. Then renormalize the remaining  $BD$  weights to sum to 1.
- Generalize the above so that  $1/var(q)$  is the adjustment factor of each set  $par_i$ , and then renormalize  $BD/var(q)$ .

A variant of the second adjustment strategy, in which an adjustment factor of zero is used when  $N_p$  is zero, worked well in our limited experience with the gene data, and that is what is used in the experiments reported in this paper. More sophisticated adjustments tied to Bayes risk are the subject of current research.

#### 4. Search

The research presented in the following sections explores two alternative methods of building the type of Bayesian classifier described in the previous sections.

The first method utilizes minimal prior knowledge regarding good parent sets for the class label and, within the Bayesian net framework, performs a simple greedy search over the entire set of genes to construct  $\mathcal{PS}$  good parent sets. The second method utilizes gene selection external to the Bayesian net framework to produce a small set  $S$  of “good genes” (like the *informative genes* of Ben-Dor *et al.* (2000) and Ben-Dor *et al.* (2001)), and then,

<sup>5</sup>Though peaked distributions which fit a large sample well tend to have better scores than dispersed distributions that fit small samples well.

within the Bayesian net framework, performs an exhaustive search of this set to find the best  $\mathcal{PS}$  subsets of  $S$  (each subset up to a specified maximum cardinality  $\mathcal{M}$ ).

#### 4.1. Minimal-Knowledge Greedy Building Methods

This family of methods ignores essentially all prior knowledge, including, in the experiments reported here, prior knowledge of which genes are “control” or “housekeeping” genes, which expression values are deemed reliable (in particular, as indicated by the  $P$ ,  $M$ , and  $A$  values in Affymetrix data), and biologically known relationships between genes. We do utilize a biological “prior” that deems it likely that only a small number of genes is necessary to classify the cases, that is, that only a small number of genes is required to render the class label conditionally independent of the remaining genes. This biological prior is necessary for any frequency-based classification method to go forward, due to sample size issues, and makes both the greedy and exhaustive searches computationally feasible. This prior is in fact supported by experiments with the current data sets in which performance—both  $BD$  and our actual classification rates—begins to diminish after a cardinality of roughly  $> 6$ . This is not quite conclusive proof, as improvement might follow disimprovement (e.g., as is exploited by simulated annealing), but this seems unlikely, especially in light of sample size issues (e.g., statistically meaningful numbers of observations of any combination of more than six gene’s expression levels is unlikely).

The version of greedy employed here proceeds in the following manner. On a designated training set (see details of the methodology in Section 5.1):

1. Use some algorithm to bin the gene expression data.
2. Determine a number  $\mathcal{K}$  of seeds, a number  $\mathcal{PS}$  of parent sets, and a maximum cardinality  $\mathcal{M}$  for the parent sets.
3. Select  $\mathcal{K}$  seed genes, based on some “goodness” criterion.
4. For each seed gene  $g_{seed}$ ,
  - a. Initialize the parent set to the singleton set  $\{g_{seed}\}$ .
  - b. Iteratively build the set to cardinality  $\mathcal{M}$  by adding one gene  $g$  at a time, chosen from the universe of all genes to maximize the  $BD$  score of  $\{\text{current set}\} \cup \{g\}$ . Maintain a list of the best  $\mathcal{PS}$  parent sets evaluated so far.
5. Construct a  $\mathcal{PS}$ -parent-set Bayesian net classifier from the list of selected parent sets as described in Section 3.6.

In Section 5.1, we specify the binning and seed selection methods used in the experiments reported in this paper.

Note that every set the greedy method evaluates, starting from each of its seeds, is a candidate for ultimate selection as one of the  $\mathcal{PS}$  parent sets—even those sets of smaller than the maximum cardinality  $\mathcal{M}$ . In particular, at every iteration, in going from cardinality  $c$  to  $c + 1$ , every extension of the best parent set of cardinality  $c$  gets a chance to be on the list of top parent sets. Consequently, some seeds may contribute more than one parent set; others may not contribute any parent sets at all.

This simple greedy method was implemented initially as a proof of concept; we suspected it would have many flaws and that we would soon replace it with more sophisticated search methods. However, it performed surprisingly well, as is attested to by both the  $BD$  scores of the best sets it finds and by the performance on our cross validation tests of the classifiers it produced (see the results in Section 5). This is not to say that avenues of potential improvements are not apparent. For example, there often is a great deal of overlap in the membership of the parent sets produced. Two or three genes tend to be present in a large fraction of the  $\mathcal{PS}$  parent sets selected. This is not necessarily a problem, but it might indicate that a nonrepresentative subspace of the set of all possible

parent sets is being searched. As is discussed in Sections 5.2 and 5.4, this effect could explain why a relatively small number of high quality parent sets are found by the algorithm.

The issue of (non)representativeness of the subspace searched is best illuminated by again considering the expression (2) that our methods are approximating. It is possible, at least in principle, that, rather than finding a representative selection of the highly weighted terms of (2), we are finding the highly weighted terms in only one region of the space. This would be problematic in the context of classification if, in other regions, the high-scoring sets contributed distributions that pushed the overall classification in an opposite direction, or, as appears to be observed in some of our testing, if not enough good parent sets are found to provide adequate coverage of the query cases. Consequently, we wish to experiment with methods that permit a wider range of regions to contribute. As an extreme, we experimented with a greedy variant that was prohibited from selecting a gene for inclusion in a second set once it was used previously. This proved too restrictive and sensitive to orderings, but several more sophisticated variants of greedy, as well as other combinatorial search heuristics, remain to be tried.

Another approach to search would mimic classical integral approximation techniques (Gander and Gautschi, 2000). In a similar learning context (Helman and Bhangoo, 1997; Helman and Gore, 1998), we employ with some success a Monte Carlo sampling method to approximate an integral representing Bayes risk. Such methods are designed to approximate an integral by sampling from regions in proportion to the amount of density contained in the region and may be adaptable to the current approximation problem.

## 4.2. External Gene Selection Methods

A second family of methods utilizes gene selection algorithms that have been developed in other contexts. This is both a promising approach to the classification problem and is indicative of how the Bayesian framework can be used to incorporate expert prior knowledge of a variety of types. As is the case with the minimal-knowledge greedy methods, we currently do not utilize prior domain knowledge about the genes; such information may, however, be discovered by our external gene selection and normalization methods and then incorporated into the framework in the form of gene selections, normalization, and binning.

The objective of external gene selection is to identify a small set of genes from which good parent sets can be constructed within a Bayesian net search procedure. By severely limiting *a priori* the size of the universe of genes to be searched for good parent sets and the maximum cardinality of the resulting parent sets, an exhaustive search for the  $\mathcal{PS}$  best parent sets (under the  $BD$  metric) can feasibly be performed. Thus, whereas the greedy method described in the previous section heuristically builds  $\mathcal{PS}$  good subsets of the universe of all genes, the external method finds the  $\mathcal{PS}$  best subsets of an intelligently restricted universe of genes.

Intuitively, good genes will be those genes whose expression values are strong indicators of a case’s classification, and we are studying a number of different methods for making such selections. The results reported in this paper are based on a strategy that computes a *separation quality* value for each gene—similar to the  $TNoM$  score described in Ben-Dor *et al.* (2000) and Ben-Dor *et al.* (2001)—and orders the genes accordingly. We then, for example, can select the genes that are the best separators.

Let  $E_1, E_2, \dots, E_n$  be the expression values for a given gene across the  $n$  cases of a training set, and let  $L_1, L_2, \dots, L_n$  be the corresponding class labels. Without loss of generality, we assume that the expression values are ordered  $E_1 \leq E_2 \leq \dots \leq E_n$  so that  $L_i$  is the class label of the  $i^{th}$  smallest expression value. The separation quality value of a gene is intended to indicate to what extent identical class labels are grouped together in  $L_1, L_2, \dots, L_n$  as a consequence of the ordering of the  $E_i$  values. Separation is considered to be perfect, for example, if the  $L_i$  labels are completely “sorted”.

Under the assumption that there are exactly two class labels,  $A$  and  $B$ , we compute separation quality as follows. Let  $Acount(i)$  be the number of  $A$  labels in  $L_1, L_2, \dots, L_i$ , and let  $Bcount(i)$  be the number of  $B$  labels in  $L_1, L_2, \dots, L_i$ . For each position  $0 \leq i \leq n$ , we can quantify the relative separation of the class labels if we were to split into the two sets  $L_1, L_2, \dots, L_i$  and  $L_{i+1}, L_{i+2}, \dots, L_n$ :

$$Separation(i) = \left| \frac{Account(i)}{Account(n)} - \frac{Bcount(i)}{Bcount(n)} \right|$$

We then define separation quality to be the best of these values:

$$SeparationQuality = \max_{1 \leq i \leq n} Separation(i)$$

Genes can be ordered by their *SeparationQuality* values, so we can talk about the  $k$  best or the  $k$  worst separators. The computed values have the following properties.

- $Account(0) = Bcount(0) = 0$
- $Bcount(i) = i - Account(i)$ , for  $0 \leq i \leq n$
- $Separation(0) = Separation(n) = 0$
- $SeparationQuality = 1$  indicates perfect separation.
- $SeparationQuality$  necessarily is  $> 0$ , since  $Separation(1)$  is  $1/Account(n)$  or  $1/Bcount(n)$ , depending on whether  $L_1$  is  $A$  or  $B$ , and we take the maximum of the *Separation* values.
- We get the same *SeparationQuality* value if we define *Account* and *Bcount* in terms of  $L_{i+1}, L_{i+2}, \dots, L_n$  instead of  $L_1, L_2, \dots, L_i$ .

We note that if the gene expression values are not distinct, then the ordering of  $E_i$  values is not unique, and the computed separation quality value will depend on the procedure used to break ties. We are considering a number of ways to pin down the ordering in the case of ties—specifically, to determine an appropriate separation quality value. We currently break these ties arbitrarily.

In addition to computing a separation quality value, we can use the same computation to propose a binning of each gene's expression values into two bins. Let  $max$  be the  $i$  value that maximizes  $Separation(i)$ , and compute

$$BinValue = \frac{E_{max} + E_{max+1}}{2},$$

which is a gene expression value that lies between the separated  $E_i$  values in the best separation. The computed *BinValue* can be used as a boundary between bins.

We note that the maximizing  $i$  value is not necessarily unique, even if the  $E_i$  values are distinct; we currently break these ties arbitrarily. We also note that  $L_{max}$  and  $L_{max+1}$  necessarily are different labels; otherwise, *SeparationQuality* could be increased by increasing or decreasing  $max$  by 1.

### 4.3 Preprocessing the Data (Normalization)

One of the advantages of the Bayesian approach is that it provides a natural mechanism to account for special domain knowledge in the construction of a classifier. Nevertheless, in our first round of experiments, we are focusing on the gene expression data, making use of minimal prior knowledge. One of the issues we are addressing in this simplified context is the preprocessing (normalization) of gene expression data before the application of our classification procedures. Because of variabilities in gene expression measurements and uncertainties about the processing done by the tools used to generate the data,<sup>6</sup> we decided to include the effect of normalization as part of our studies.

<sup>6</sup>Affymetrix Microarray Suite (MAS) Version 4.0.

Our approach to normalization is to consider, for each case, the average expression value over some designated set of genes, and to scale each case so that this average value is the same for all cases. This approach allows our analysis to concentrate on relative gene expression values within a case by standardizing a reference point between cases. For example, if the expression value within a case of certain genes  $g_i$  relative to the expression value of some reference point gene  $g$  is an effective class discriminator, then it suffices simply to consider these  $g_i$  values, provided cases have first been normalized to a common  $g$  value. The key difference between the normalization strategies we considered is the choice of the reference point gene  $g$ , or, more generally, the choice of a set  $R$  of reference point genes. While selecting an appropriate set  $R$  could provide a good opportunity to take advantage of special knowledge of the underlying domain, consistent with our desire to focus first on raw data in the absence of prior knowledge, we use here a simple selection method based on the *SeparationQuality* value already discussed. In particular, we set  $R$  to be the  $k$  worst separators—that is, genes with the lowest *SeparationQuality* values—for some number  $k$ . The motivation for this choice of  $R$  is that, as our experiments indicate, a suitable reference point can be found as the average of the expression values of genes that are independent of the class label for which we are trying to develop a classifier. Further, normalizing with respect to such genes will not discard information that might be valuable in class discrimination. Choosing the  $k$  worst separators for normalization is a heuristic for identifying genes likely to be independent of the class label.

In summary, the normalization algorithm we used is as follows.

1. Let  $R$  consist of the  $k$  worst separator genes, as described above.
2. Let  $A$  represent the target average value for the genes in  $R$ ;  $A$  may be chosen arbitrarily, since its value does not affect any aspects of the computation.
3. For each case  $C$ ,
  - a. Compute the average value,  $Ave_{R,C}$ , of the expression values in case  $C$  for the genes in  $R$ .
  - b. Multiply every expression value of case  $C$  by the scaling factor  $A/Ave_{R,C}$ .

We took  $k$  to be a parameter to be learned in the course of training and experimented with several different values accordingly. The results of these experiments against training data are reported in Section 5.3; Section 5.4 reports how well a choice of  $k$  made against training data generalizes to an out-of-sample test set.

## 5. Results

The MIT leukemia data (Golub *et al.*, 1999) and the Princeton colon cancer data (Alon *et al.*, 1999) are considered. The MIT data consists of 7,129 gene expression values per case. The Princeton data is provided in a heavily pruned form, consisting of only 2,000 genes per case.

### 5.1. Experimental Methodology

In order to avoid any possibility of overfitting our results to specific data sets, we set aside from each data set a fraction of cases, forming a *test set*. For the MIT data set, a partition into 38 training cases and 34 test cases (our set aside, out-of-sample cases) is provided on the MIT Web site. The Princeton Web site provides a single data set of 62 cases. We randomly partitioned this data set into 38 training cases and 24 set aside test cases. The test sets were not examined *at all* in the course of algorithm development, nor to establish parameter settings, and were considered only after our best methods, along with their parameter settings, were identified through a cross validation methodology (detailed below) on the training sets. Results of our best method—as identified against the training sets only—run against the set aside test sets are reported in Section 5.4.

We now describe the cross validation methodology that was applied to the 38-case training sets in order to develop our methods and to indicate which techniques would be the most promising to pursue. In particular,

our initial evaluation of a classifier building method under development employed “leave one out” (*LOO*) cross validation. On each experiment, a method would train on 37 cases, building a classifier to be used to classify the single left out query case; the build/evaluate cycle is repeated 38 times, once for each “fold” created by leaving out of the training a different query case.

Care must be taken during development that the methods used in the classifier construction process not exploit *any* knowledge of the left out query case it is to be evaluated on. That is, any method applied to build the classifier must be applicable when we turn attention to the set aside test set (or to an actual set of query cases for which a classification is desired), at which time knowledge of the query’s class label, of course, is unavailable.

This requirement implies, for example:

- Gene selection by external means must be repeated on each of the 38 folds, without being exposed to the left out case to be used as a query in the evaluation.
- Similarly, if normalization or binning is to use label knowledge, it must not be exposed to the left out case, and hence must be repeated for each fold. If, however, a binning algorithm does not use knowledge of labels (as is the case of the algorithm used in connection with the greedy construction), it may inspect the entire training set, since in an actual classification application, the binning algorithm could inspect the non-label fields (genes) of the cases to be classified at the time these cases are presented for analysis.

### **Greedy Parent Set Construction**

The LOO cross validation setup for the greedy method takes the following form:

1. Let  $T$  represent the full training set (e.g., of 38 cases).
2. Bin  $T$ , without using label knowledge.
3. For each  $q_i \in T$ , define fold  $F_i = (T - \{q_i\})$ ,
  - a. Select  $\mathcal{K}$  seeds against  $F_i$ .
  - b. Use the greedy method to construct  $\mathcal{PS}$  good sets (under  $BD$ ) up to cardinality  $\mathcal{M}$  against  $F_i$ , starting from each seed.
  - c. Compute the variance in each set’s induced distribution of  $q_i$ ’s unknown label, and adjust the  $BD$  score of each set to form a  $\mathcal{PS}$ -set classifier.
  - d. Classify  $q_i[C]$  as the most likely value, given  $q_i[genes]$  under the classifier’s distribution.
  - e. Compute the error and uncertainty in the classification for fold  $F_i$ .
4. Report the average error and uncertainty rates across the folds.

The information reported in Step 4 is derived from the constructed classifiers’ induced distributions. In particular, the classifier constructed for each fold  $F_i$  specifies a conditional posterior distribution  $Pr\{q[C] = c_k \mid q[genes] = \langle e \rangle\}$  for a query case’s class label. In the current experiments, the class label is binary, and  $q$  is classified as belonging to the class with higher posterior; if value = 0.5, no classification is possible. An error occurs if  $q[C]$  is the lower probability class.

Uncertainty is a measure of the strength of a classification. If  $Pr\{q[C] = c_k \mid q[genes] = \langle e \rangle\}$  is near 1.0, the classification is strong, whereas if it is near 0.5, it is weak. On each fold, we compute the “probability of error” as well as the 0/1 misclassification indicator. In particular, probability of error is given by  $(1.0 - (\text{the probability the classifier assigns to the true class } q[C]))$ .

For the experiments reported in Section 5.2 and 5.4, we utilized the following relatively simple binning (Step 2) and seed selection (Step 3.a) techniques.

Binning: As is indicated in Section 3.1, practical Bayesian net methods require a discretization of the expression values. Following most gene expression researchers, we partition values into three ranges: “under-”, “average-”, and “over-” expressed. Our partitioning method for greedy creates a tertiary binning for each gene  $g$  as

$$\begin{aligned} &(-\infty, (mean(g) - n_{low} \times \sigma(g)), \\ &[mean(g) - n_{low} \times \sigma(g), mean(g) + n_{high} \times \sigma(g)], \\ &(mean(g) + n_{high} \times \sigma(g), \infty), \end{aligned}$$

where the mean  $mean(g)$  and standard deviation  $\sigma(g)$  of each gene’s  $g$  expression values are computed over all cases. The choices of  $n_{low}$  and  $n_{high}$  are made through experimentation on the training data. Once selected, these are fixed and used without modification on the set aside test data; otherwise, we would run the risk of overfitting to the data. For the MIT data, setting  $n_{low} = n_{high} = 1.0$  worked well, and there was little sensitivity in the cross validation results. In the Princeton data, there was far more sensitivity in the cross validation, and a limited search arrived at the settings  $n_{low} = 1.25$  and  $n_{high} = 0.4$ .

Seed selection: Singleton parent sets  $\{g\}$  are formed for each gene  $g$  and the  $BD$  score obtained. The genes corresponding to the  $\mathcal{K}$  highest scoring parent sets are used as seeds.

### External Gene Selection Plus Exhaustive Parent Set Construction

The LOO cross validation setup for external gene selection takes the following form:

1. Let  $T$  represent the full training set (e.g., of 38 cases).
2. For each fold defined by  $F_i = (T - \{q_i\})$ ,
  - a. Use an external method against  $F_i$  to normalize expression values and select a set  $S$  of  $\mathcal{N}$  genes.
  - b. Bin  $F_i$ , possibly using information returned by gene selection.
  - c. Exhaustively search the set  $S$  for the best  $\mathcal{PS}$  subsets (of cardinality up to  $\mathcal{M}$ ) under the  $BD$  scoring metric.
  - d. Compute the variance in each set’s induced distribution of  $q_i$ ’s unknown label, and adjust the  $BD$  score of each set to form a  $\mathcal{PS}$ -set classifier.
  - e. Classify  $q_i[C]$  as the most likely value, given  $q_i[genes]$  under the classifier’s distribution.
  - f. Compute the error and uncertainty in the classification for fold  $F_i$ .
3. Report the average error and uncertainty rates across the folds.

In our experiments, we employed the external gene selection, normalization, and binning methods described in Section 4.2. In particular, the external gene selection algorithm is invoked on each fold with the following effect:

- The algorithm normalizes the cases in  $F_i$  using the  $k$  genes with the lowest *SeparationQuality* as controls.
- The algorithm returns the  $\mathcal{N}$  genes with the highest *SeparationQuality*.
- The algorithm returns a binary bin boundary for each selected gene, corresponding to where the maximum separation value is obtained.

Once results of the external gene selection algorithm are returned for a fold, an exhaustive search is performed (on a normalized and binned  $F_i$ ) for the best  $\mathcal{PS}$  parent sets, from which the Bayesian net classifier is formed.

Note that the instantiation of the steps of either methodology with specific algorithms defines a classifier building method. When run on a specific training set (or fold of a training set), it yields a  $\mathcal{PS}$ -set classifier, which in turn yields a posterior class distribution. This distribution can then be used to classify query cases with unknown labels, assuming that the query cases are drawn from the same distribution which underlies the training set. We emphasize that it is the building method, not the particular classifiers built on a run against a training set (or fold of a training set), that is being assessed.

## 5.2. Cross Validation Results with Greedy

In tests of the greedy method, we studied the effects of varying the number  $\mathcal{PS}$  of sets used in the classifier. We held fixed at  $\mathcal{M} = 5$  the maximum cardinality and, due to computational considerations, the number of seeds at  $\mathcal{K} = 60$ .

The following two tables summarize, respectively, results with the Princeton and MIT training sets. Each row of the tables summarizes, for a fixed  $\mathcal{PS}$ , the LOO cross validation test results for the 38 cases of the respective training set. The  $qMax$  result appearing at the end of each table is discussed below.

Legend:

- $\mathcal{PS}$  : Number of parent sets used.
- $APE$  : Average probability error per fold.
- $MIS$  : Number of misclassifications.
- $ERR$  : Total error count (misclassifications + nonclassifications).
- $TER$  : Total error rate (including both misclassifications and nonclassifications).

$\mathcal{PS}$	$APE$	$MIS$	$ERR$	$TER$
1	0.184212	4	10	0.263158
5	0.169929	7	7	0.184211
10	0.259123	12	12	0.315789
20	0.312331	14	14	0.368421
60	0.329858	13	13	0.342105
300	0.340612	13	13	0.342105
500	0.346113	14	14	0.368421
$qMax$	0.289474	11	11	0.289474

**Table 1. Princeton training data** ( $n_{low} = 1.25, n_{high} = 0.4$ ).

$\mathcal{PS}$	$APE$	$MIS$	$ERR$	$TER$
1	0.315791	0	24	0.631579
5	0.193975	1	14	0.368421
10	0.140994	1	9	0.236842
20	0.067464	2	3	0.078947
60	0.070245	3	3	0.078947
300	0.089030	3	3	0.078947
500	0.118584	5	5	0.131579
$qMax$	0.157897	6	6	0.157897

**Table 2. MIT training data** ( $n_{low} = 1.0, n_{high} = 1.0$ ).

The tables indicate an initial increase in quality as  $\mathcal{PS}$  increases, then a leveling off and ultimate decrease in quality. The most interesting result is the significant increase in quality over just a single set ( $\mathcal{PS} = 1$ , the *maximum a posteriori solution*), which is a prevalent Bayesian net methodology for learning distributions. As predicted from the discussion in Section 3.3, a single parent set does not provide adequate coverage of gene expression combinations in the query case, leading to a large number of non classifications.

To establish that the polling effect noted in Section 3.3 is real and significant, we also conducted experiments labeled “*qMax*”. Here, 500 sets are built as with  $\mathcal{PS} = 500$ , but for each query case  $q$ , the single parent set with the highest variance adjusted score is used to classify  $q$ . Note that this query-specific set selection from the 500 always selects (if available, which is the case in all our cross validation runs) a set in which  $q$ ’s combination of expression values appears in the training set, eliminating the no-classification errors. That this method underperforms the best  $\mathcal{PS} > 1$  methods indicates that the blending of distributions contributes to the quality of the classification. Examination of the details of the computations performed by the classifier also indicates that, in many cases, the distributions induced by the parent sets exert competing effects on the classification, and that the weighting resolution generally leads to a correct classification.

We speculate that the degradation in classification quality for  $\mathcal{PS}$  above a threshold is caused by the potentially unrepresentative search performed by our simple greedy algorithm, as alluded to in Section 4.1—greedy, being unable to construct enough high scoring sets, must “fill” the classifier with many low scoring (and, hence, worse fitting to the observational data) sets which contribute inaccurate distributions. This explanation is supported by the near monotonic increase in quality reported in Section 5.3 for the exhaustive search following external gene selection. This suggests that refinements to greedy as proposed in Section 4.1 could well obtain overall improvements, especially as is noted in Section 5.4 when we discuss the results of the greedy-built classifiers against the out-of-sample test set.

### 5.3. Cross Validation Results with External Gene Selection

In tests of the external gene selection methods, we studied the effects of varying both  $\mathcal{PS}$  and the fraction  $\mathcal{W}$  of genes used as controls in normalization. As with greedy, we held fixed the maximum cardinality  $\mathcal{M}$  at 5. For computational reasons, the number of genes selected was fixed at 30.

The following two tables summarize, respectively, results with the Princeton and MIT training sets. Each row of the tables summarizes, for a fixed  $\mathcal{W}$  and  $\mathcal{PS}$ , the LOO cross validation test results for the 38 cases of the respective training set. As is the case for Tables 1 and 2, the *qMax* result at the end of each of Tables 3 and 4 is for 500 available parent sets and with  $\mathcal{W}$  set at a value which produced generally good results across the  $\mathcal{PS}$  values for the multi-set classifiers.

Legend:

- $PS$  : Number of parent sets used.
- $\mathcal{W}$  : Fraction of genes used as controls for normalization.
- $APE$  : Average probability error per fold.
- $MIS$  : Number of misclassifications.
- $ERR$  : Total error count (misclassifications + nonclassifications).
- $TER$  : Total error rate (including both misclassifications and nonclassifications).

$PS$	$\mathcal{W}$	$APE$	$MIS$	$ERR$	$TER$
1	0.000000	0.394735	15	15	0.394737
1	0.100000	0.480700	17	20	0.526316
1	0.250000	0.328950	8	17	0.447368
1	0.400000	0.302636	8	15	0.394737
1	0.550000	0.328950	7	18	0.473684
1	0.700000	0.263162	7	13	0.342105
1	0.850000	0.499997	18	20	0.526316
1	1.000000	0.499994	16	22	0.578947
5	0.000000	0.393831	14	14	0.368421
5	0.100000	0.376276	14	14	0.368421
5	0.250000	0.287669	9	11	0.289474
5	0.400000	0.267520	9	11	0.289474
5	0.550000	0.241729	9	10	0.263158
5	0.700000	0.261501	9	10	0.263158
5	0.850000	0.455024	17	17	0.447368
5	1.000000	0.333537	10	13	0.342105
10	0.000000	0.377660	15	15	0.394737
10	0.100000	0.398858	15	15	0.394737
10	0.250000	0.334334	12	12	0.315789
10	0.400000	0.261875	9	11	0.289474
10	0.550000	0.221307	8	9	0.236842
10	0.700000	0.270484	9	9	0.236842
10	0.850000	0.410469	14	14	0.368421
10	1.000000	0.303383	10	11	0.289474

**Table 3. Princeton training data.**

$\mathcal{P}\mathcal{S}$	$\mathcal{W}$	$APE$	$MIS$	$ERR$	$TER$
20	0.000000	0.377660	15	15	0.394737
20	0.100000	0.402184	16	16	0.421053
20	0.250000	0.302113	11	11	0.289474
20	0.400000	0.251675	9	9	0.236842
20	0.550000	0.215504	7	8	0.210526
20	0.700000	0.265321	9	9	0.236842
20	0.850000	0.361076	12	12	0.315789
20	1.000000	0.325153	11	12	0.315789
60	0.000000	0.350131	12	12	0.315789
60	0.100000	0.375262	14	14	0.368421
60	0.250000	0.290695	10	10	0.263158
60	0.400000	0.233612	9	9	0.236842
60	0.550000	0.204675	7	7	0.184211
60	0.700000	0.249359	8	8	0.210526
60	0.850000	0.358279	12	12	0.315789
60	1.000000	0.286617	10	11	0.289474
300	0.000000	0.344514	13	13	0.342105
300	0.100000	0.358541	14	14	0.368421
300	0.250000	0.297478	11	11	0.289474
300	0.400000	0.223621	7	7	0.184211
300	0.550000	0.204802	7	7	0.184211
300	0.700000	0.237995	8	8	0.210526
300	0.850000	0.317356	12	12	0.315789
300	1.000000	0.249347	9	9	0.236842
500	0.000000	0.341484	13	13	0.342105
500	0.100000	0.351571	14	14	0.368421
500	0.250000	0.293802	12	12	0.315789
500	0.400000	0.218802	7	7	0.184211
500	0.550000	0.206535	6	6	0.157895
500	0.700000	0.231278	8	8	0.210526
500	0.850000	0.301052	11	11	0.289474
500	1.000000	0.251559	9	9	0.236842
$qMax$	0.550000	0.210529	8	8	0.210526

**Table 3. Princeton training data (continued).**

$\mathcal{P}\mathcal{S}$	$\mathcal{W}$	$APE$	$MIS$	$ERR$	$TER$
1	0.000000	0.065801	1	4	0.105263
1	0.100000	0.052644	1	3	0.078947
1	0.250000	0.065801	1	4	0.105263
1	0.400000	0.078959	1	5	0.131579
1	0.550000	0.078959	1	5	0.131579
1	0.700000	0.078959	1	5	0.131579
1	0.850000	0.065802	1	4	0.105263
1	1.000000	0.078959	1	5	0.131579
5	0.000000	0.072555	3	3	0.078947
5	0.100000	0.053353	2	2	0.052632
5	0.250000	0.072555	3	3	0.078947
5	0.400000	0.080379	3	3	0.078947
5	0.550000	0.080379	3	3	0.078947
5	0.700000	0.080379	3	3	0.078947
5	0.850000	0.061176	2	2	0.052632
5	1.000000	0.080378	3	3	0.078947
10	0.000000	0.072554	3	3	0.078947
10	0.100000	0.053351	2	2	0.052632
10	0.250000	0.072554	3	3	0.078947
10	0.400000	0.080378	3	3	0.078947
10	0.550000	0.080378	3	3	0.078947
10	0.700000	0.080378	3	3	0.078947
10	0.850000	0.061175	2	2	0.052632
10	1.000000	0.083038	3	3	0.078947

**Table 4. MIT training data.**

$\mathcal{P}\mathcal{S}$	$\mathcal{W}$	$APE$	$MIS$	$ERR$	$TER$
20	0.000000	0.072553	3	3	0.078947
20	0.100000	0.053351	2	2	0.052632
20	0.250000	0.072553	3	3	0.078947
20	0.400000	0.080377	3	3	0.078947
20	0.550000	0.080377	3	3	0.078947
20	0.700000	0.080377	3	3	0.078947
20	0.850000	0.061174	2	2	0.052632
20	1.000000	0.084275	3	3	0.078947
60	0.000000	0.070544	3	3	0.078947
60	0.100000	0.051839	2	2	0.052632
60	0.250000	0.071990	3	3	0.078947
60	0.400000	0.069324	3	3	0.078947
60	0.550000	0.070813	3	3	0.078947
60	0.700000	0.070774	3	3	0.078947
60	0.850000	0.050437	2	2	0.052632
60	1.000000	0.069833	3	3	0.078947
300	0.000000	0.057444	2	2	0.052632
300	0.100000	0.059049	2	2	0.052632
300	0.250000	0.072465	3	3	0.078947
300	0.400000	0.074483	3	3	0.078947
300	0.550000	0.074229	3	3	0.078947
300	0.700000	0.075196	3	3	0.078947
300	0.850000	0.056310	2	2	0.052632
300	1.000000	0.050879	2	2	0.052632
500	0.000000	0.065868	2	2	0.052632
500	0.100000	0.068942	2	2	0.052632
500	0.250000	0.080150	3	3	0.078947
500	0.400000	0.079164	3	3	0.078947
500	0.550000	0.078501	3	3	0.078947
500	0.700000	0.078683	3	3	0.078947
500	0.850000	0.074403	2	2	0.052632
500	1.000000	0.068241	2	2	0.052632
$qMax$	0.100000	0.052644	2	2	0.052632

**Table 4. MIT training data (continued).**

Unlike the case for greedy selection, the results of Tables 3 and 4 demonstrate that there is a steady improvement for the Princeton data as  $\mathcal{PS}$  increases, and near flat behavior for the MIT data for  $\mathcal{PS} \geq 60$ . Again, the  $qMax$  experiments (for the Princeton data) and inspection of the detailed results provide further evidence that the blending provided by a large number of parent sets has a positive impact on classifier quality.

The tables indicate different best values across the two training sets for the fraction  $\mathcal{W}$  of control genes used in expression-level normalization, and a greater sensitivity to this value in the Princeton training data. This may be indicative of differences in experimental conditions, analysis preprocessing, and so forth. That we can, without the benefit of descriptive procedural information as input, discover through methodical application of cross validation good normalization parameters for each data set is a significant finding. The results against the test set presented in the following section indicate that these findings are not simply an overfitting to the training data, but truly a learning of the underlying processes that generalizes well.

#### 5.4. Out-of-Sample Test Set Results

Only after running the above experiments on the training sets did we turn attention to the test sets. Our primary interest is to select the *single method* which performed best (lowest total error rate,  $TER$ ) in the cross validation experiments and assess its classification rate on the out-of-sample test sets. In this way, we avoid a “selection effect” in which one of several methods run against the test set performs well.

Inspection of the tables of Sections 5.2 and 5.3 identifies the external gene selection method as being preferable to the minimal knowledge greedy method in building parent sets for the Bayesian net classifier. Since we have data from two different experimental contexts, it is proper to select the parameters for the selected method (i.e.,  $\mathcal{PS}$  and  $\mathcal{W}$ ) based on performance in the cross validation trials on each training set; such parameter setting would of course be performed in an actual classification application in which we had access to training, but not query, cases in advance.

#### External Gene-Selection Method Against Test Data

Inspection of the tables in Section 5.3 indicates that, against the Princeton training set, the best setting is  $\mathcal{PS} = 500$  (number of parent sets to be used in the Bayesian net classifier) and  $\mathcal{W}=0.55$  (control list fraction for normalization). Against the MIT training set, several parameter settings resulted in the minimal  $TER$  of 0.052632. Somewhat arbitrarily, we selected  $\mathcal{PS}=300$  and  $\mathcal{W} = 0.85$ .<sup>7</sup> Using *only these settings*, we built the classifiers by training against the 38 cases of each of the two training sets and used the resulting classifiers to classify the cases of the respective test sets.

The results are exhibited in Table 5 and are extremely good. The classifier had nearly identical error rates against the MIT training and test sets (0.05 for training versus 0.06 for test) and a significantly lower error rate against the Princeton test set (0.16 for training versus 0.08 for test). The results strongly suggest that our multi-parent-set Bayesian net classifiers employing external gene selection and normalization algorithms are able to learn from training data underlying distributions which generalize extremely well to out-of-sample query cases whose classifications are of biological and clinical significance.

Test Set	Cases	$APE$	$MIS$	$ERR$	$TER$
Princeton	24	0.142092	2	2	0.083333
MIT	34	0.085831	2	2	0.058824

**Table 5. Out-of-sample results with external gene selection.**

<sup>7</sup>While we chose our single run to be made against the test set with  $\mathcal{PS} = 300$  and  $\mathcal{W} = 0.85$ , in order to assess the sensitivity of the results to this somewhat arbitrary choice of settings from among settings achieving equally good  $TER$ , we later ran against the test set with several other settings which achieved the same  $TER$  against the training data. The majority of those settings tried also incurred the same number 2 of misclassification errors as those reported here, while a few others incurred 3 misclassifications errors.

### Minimal-Knowledge Greedy Methods Against Test Data

After obtaining the results reported in the previous subsection for the external methods, we decided also to run our greedy methods against the test sets. Since the greedy method’s results in the cross validation experiments were almost as good as the external gene selection methods, we consider this to be an interesting avenue of research as well. We report the results here in order to indicate potential directions for future work.

Table 6 reports the results against the two test sets of Bayesian net classification using the greedy construction method. The only parameter considered in the cross validation against the training set was  $\mathcal{PS}$ , with the best settings found to be  $\mathcal{PS} = 20$  for the MIT training set and  $\mathcal{PS} = 5$  for the Princeton training set.

Test Set	Cases	<i>APE</i>	<i>MIS</i>	<i>ERR</i>	<i>TER</i>
Princeton	24	0.145834	1	6	0.250000
MIT	34	0.279412	7	12	0.352941

**Table 6. Out-of-sample results with greedy selection.**

Against the Princeton test set, the error rate was similar to the rate against the training set (0.18 for training versus 0.25 for test), but it was significantly higher against the MIT test set (0.08 for training versus 0.35 for test). We speculate that two sources of this lack of generalization, especially in the MIT data, are our failure to normalize the data for the greedy experiments and the use of an overly rigid binning method. This conjecture is consistent with the high number of “nonclassifications” against the test sets. Note also that the MIT data was provided as two distinct data sets. Procedural differences in experimental preparation and processing of the output between the sets (Golub *et al.*, 1999) may have hampered the greedy method because it fails to normalize across the sets. In the case of the Princeton data, where a single data set is randomly split, performance against the test set was much more comparable to that of the training set.

Consequently, one avenue of future research is to include in the greedy method a normalization procedure similar to that employed by the external gene selection method. Also, as noted in Section 4.1, there is a concern that the greedy search may not provide a good representation of the space of possible parent sets. We speculated that this might be the cause of the degradation observed in the cross validation experiments for large values of  $\mathcal{PS}$ . Note that the exhaustive (and, hence, completely representative) search of the universe of externally selected genes resulted in large  $\mathcal{PS}$ s performing best. The greedy method’s use of small values of  $\mathcal{PS}$ , in combination with the failure to normalize, certainly contributes to the large number of non-classifications in the test set. Hence, modifying the search to be more representative, as discussed in Section 4.1, potentially could give minimal-knowledge searches such as greedy access to more good parent sets, thereby addressing the large number of failure-to-classify errors that were observed.

## 6. Summary and Future Work

We have presented a methodology for applying Bayesian nets to the problem of classifying clinical cases from their gene expression profiles. While Bayesian nets have been applied previously to identify relationships among genes and have been proposed as classifiers for other problem domains, we have outlined new methods for classification particularly well suited to gene expression data. Through a systematic experimental design, we demonstrated that these classifiers, trained by means of a cross-validation methodology, generalize extremely well to out-of-sample test data. In particular, we achieved error rates of 92% and 94% on out-of-sample partitions of the MIT leukemia and Princeton colon cancer data sets, respectively.

Our Bayesian net classifiers are built by constructing alternative parent sets for the class label node and use a posterior probability and variance-weighted blending of the resulting distributions. This blending of the distributions induced by the competing hypotheses embodied by the alternative parent sets was seen in our experimental results to yield improvements over the so called *maximum a posteriori solution*, in which only the single most

likely hypothesis is used. We experimented with two methods for searching for good parent sets: a simple greedy search of the universe of all genes and an exhaustive search of a universe of genes selected by a separation heuristic. The latter method produced better performing parent sets in the experiments reported here. This method also employs a novel expression-level normalization scheme based on algorithmically discovered control genes. Current work is considering improvements to both methods for parent set construction and to normalization. We are exploring also how other aspects of the problem—value binning and gene clustering, for example—can be studied within the framework.

We believe that Bayesian approaches to gene expression analysis, such as those described here and in Friedman *et al.* (1999), Friedman *et al.* (2000) and Pe'er *et al.* (2001), have enormous potential, not simply because of the quality of the results achieved so far, but also because the mathematically-grounded formalism provides the opportunity to expand systematically the range of problems treated, integrating newly developed algorithmic techniques with an ever-increasing base of domain knowledge. Thus, results such as those reported here, while significant in their own right, are only the first steps toward the ultimate construction of rigorous and comprehensive models that promise to be of great scientific and clinical import.

### Acknowledgments

This work has been supported by grants from the D.H.H.S. National Institutes of Health/National Cancer Institute (CA88361), the W.M. Keck Foundation, and National Tobacco Settlement funds to the State of New Mexico Provided to UNM for Genomics and Bioinformatics. We are grateful for the computational support that has been provided by the UNM High Performance Computing Education and Research Center (HPCERC).

### References

- Alizadeh, A., Eisen, M., Davis, R., Ma, C., Lossos, I., Rosenwald, A., Boldrick, J., Sabet, H., Tran, T., Yu, X., Powell, J., Yang, L., Marti, G., Moore, T., Hudson, J. (Jr), Lu, L., Lewis, D., Tibshirani, R., Sherlock, G., Chan, W., Greiner, T., Weisenburger, D., Armitage, J., Warnke, R., Levy, R., Wilson, W., Grever, M., Byrd, J., Botstein, D., Brown, P., and Staudt, L. 2000. Distinct types of diffuse large B-cell lymphoma identified by gene expression profiling. *Nature* 403, 503–511.
- 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 *Proc. National Academy of Sciences USA* 96(12), 6745–6750.
- Bhattacharjee, A., Richards, W., Staunton, J., Li, C., Monti, S., Vasa, P., Ladd, C., Beheshti, J., Bueno, R., Gillette, M., Loda, M., Weber, G., Mark, E., Lander, E., Wong, W., Johnson, B., Golub, T., Sugarbaker, D., and Meyerson, M. 2001. Classification of human lung carcinomas by mRNA expression profiling reveals distinct adenocarcinoma subclasses. In *Proc. National Academy of Sciences* 98(24), 13790–13795.
- Ben-Dor, A., Bruhn, B., Friedman, N., Nachman, I., Schummer, M., and Yakhini, Z. 2000. Tissue classification with gene expression profiles. *J. Computational Biology* 7(3,4), 559–584.
- Ben-Dor, A., Friedman, N., and Yakhini, Z. 2001. Class discovery in gene expression data. In *Proc. Fifth Annual International Conference on Computational Biology*, 31–38, ACM Press.
- Brown, P., and Botstein, D. 1999. Exploring the new world of the genome with DNA microarrays. *Nature Genetics* 21, 33–37.
- Buntine, W. 1991. Theory refinement for Bayesian networks. In *Proc. Seventh Conference on Uncertainty in Artificial Intelligence*, 52–60, Morgan Kaufmann.
- Buntine, W. 1996. A guide to the literature on learning probabilistic networks from data. *IEEE Trans. on Knowledge and Data Engineering* 8, 195–210.

- Cooper, G., and Herskovita, E. 1992. A Bayesian method for the induction of probabilistic networks from data. *Machine Learning* 9, 309–347.
- Cozman, F. 2001. <http://www-2.cs.cmu.edu/~javabayes/Home/>.
- D’haeseleer, P. 2000. Reconstructing gene networks from large scale gene expression data. Ph.D. dissertation. Computer Science Department, University of New Mexico. <http://www.cs.unm.edu/~patrik>.
- Dawid, A.. 1992. Applications of a general propagation algorithm for a probabilistic expert system. *Statistics and Computing* 2, 25–36.
- Dawid, A., and Lauritzen, S. 1993. Hyper Markov laws in the statistical analysis of decomposable graphical models. *Annals of Statistics* 21, 1272–1317.
- Decter, R. 1996. Bucket elimination: a unifying framework for probabilistic inference. In *Proc. Twelfth Conference on Uncertainty in Artificial Intelligence*, 211–219, Morgan Kaufmann.
- DeRisi, J., Iyer, V., and Brown, P. 1997. Exploring the metabolic and genetic control of gene expression on a genomic scale. *Science* 278, 680–686.
- Duda, R., and Hart, P. 1973. *Pattern classification and scene analysis*. John Wiley and Sons, New York.
- Eisen, M., Spellman, P., Botstein, D., and Brown, P. 1998. Cluster analysis and display of genome-wide expression patterns. In *Proc. National Academy of Sciences USA* 95, 14863–14867.
- Friedman, N., and Goldszmidt, M. 1996. Learning Bayesian networks with local structure. In *Proc. Twelfth Conference on Uncertainty in Artificial Intelligence*, 211–219, Morgan Kaufmann.
- Friedman, N., Geiger, D., and Goldszmidt, M. 1997. Bayesian network classifiers. *Machine Learning* 29, 131–163.
- Friedman, N., Nachman, I., and Pe’er, D. 1999. Learning Bayesian network structure from massive datasets: the sparse candidate algorithm. In *Proc. Fifteenth Conference on Uncertainty in Artificial Intelligence*, 196–205, Morgan Kaufmann.
- Friedman, N., Linial, M., Nachman, I., and Pe’er, D. 2000. Using Bayesian networks to analyze expression data. *J. Computational Biology* 7(3,4), 601–620.
- Gander, W., and Gautschi, W. 2000. Adaptive Quadrature Revisited. *BIT* 40(1), 84–101.
- Getz, G., Levine, E., and Domany, E. 2000. Coupled two-way clustering analysis of gene microarray data. In *Proc. National Academy of Sciences* 97(22), 12079–12084.
- Golub, T., Slonim, D., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J., Coller, H., Loh, M., Downing, J., Caligiuri, M., Bloomfield, C., and Lander, E. 1999. Molecular classification of cancer: Class discovery and class prediction by gene expression monitoring. *Science* 286, 531–537.
- Heckerman, D., Geiger, D., and Chickering, D. 1995. Learning Bayesian networks: The combination of knowledge and statistical data. *Machine Learning* 20, 197–243.
- Helman, P., and Bhangoo, J. 1997. A statistically based system for prioritizing information exploration under uncertainty. *IEEE Trans. on Systems, Man, and Cybernetics* 27(4), 449–466, 1997.
- Helman, P., and Gore, R. 1998. Prioritizing information for the discovery of phenomena. *J. of Intelligent Information Systems* 11(2), 99–138.
- Ibrahim, J., Chen, M., and Gray, R. 2002. Bayesian models for gene expression with DNA microarray data. *J. American Statistical Association* 97(457), 88–99.
- Jensen, F., Lauritzen, S., and Olesen, K. 1990. Bayesian updating in causal probabilistic networks by local computations, *Computational Statistics Quarterly* 4, 269–282.

- Jensen, F. 2001. *Bayesian networks and decision graphs*. Springer-Verlag, New York.
- Khan, J., Wei, J., Ringner, M., Saal, L., Ladanyi, M., Westermann, F., Berthold, F., Schwab, M., Antonescu, C., Peterson, C., and Meltzer, P. 2001. Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks. *Nature Medicine* 7, 673–679.
- Lam, W., and Bacchus, F. 1994. Learning Bayesian belief networks: an approach based on the MDL principle, *Computational Intelligence* 10, 269–293.
- Langley, P., Iba, W., and Thompson, K. 1992. An analysis of Bayesian classifiers. In *Proc. Tenth National Conference on Artificial Intelligence*, 223–228, AAAI Press.
- Lauritzen, S., and Spiegelhalter, D. 1988. Local computations with probabilities on graphical structures and their applications to expert systems, *J. Royal Statistical Society Series B* 50, 157–224.
- Lockhart, D., Dong, H., Byrne, M., Follettie, M., Gallo, M., Chee, M., Mittmann, M., Wang, C., Kobayashi, M., Horton, H., and Brown, E. 1996. Expression monitoring by hybridization to high-density oligonucleotide arrays. *Nature Biotechnology* 14(12), 1675.
- Madsen, A., and Jensen, F. 1999. Lazy propagation: a junction tree inference algorithm based on lazy evaluation, *Artificial Intelligence* 113, 203–245.
- Murphy, K., and Mian, S. 1999. Modelling gene expression data using dynamic Bayesian networks. Technical Report, Computer Science Division, University of California, Berkeley. <http://www.cs.berkeley.edu/~murphyk/Papers/ismb99.ps.gz>.
- Pearl, J. 1988. *Probabilistic reasoning for intelligent systems*. Morgan Kaufmann, San Francisco.
- Pearl, J., and Verma, T. 1991. A theory of inferred causation. In *Knowledge Representation and Reasoning: Proc. Second International Conference*, 411–452, Morgan Kaufmann.
- Pe'er, D., Regev, A., Elidan, G., and Friedman, N. 2001. Inferring subnetworks from perturbed expression profiles. *Bioinformatics* 17, S215–S224.
- Pomeroy, S., Tamayo, P., Gassenbeek, 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., Blegel, J., Pogglo, 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. *Nature* 415, 436–442.
- Rigoutsos, I., Floratos, A., Parida, L., Gao, Y., and Platt, D. 2000. The emergence of pattern discovery techniques in computational biology. *Metabolic Engineering* 2(3), 159–177.
- Schena, M., Shalon, D., Davis, R., and Brown, P. 1995. Quantitative monitoring of gene expression patterns with a complementary DNA microarray. *Science, New Series* 270(5235), 467–470.
- Shafer, G., and Shenoy, P. 1990. Probability propagation. *Annals of Mathematics and Artificial Intelligence* 2, 327–352.
- Spiegelhalter, D., Dawid, D., Lauritzen, S., and Cowell, R. 1993. Bayesian analysis in expert systems. *Statistical Science* 8, 219–292.
- Tavazoie, S., Hughes, J., Campbell, M., Cho, R., and Church, G. 1999. Systematic determination of genetic network architecture. *Nature Genetics* 22, 281–285.
- Tobin, F., Damian-Iordache, V., and Greller, L. 1999. Towards the reconstruction of gene regulatory networks. In *Technical Proc. 1999 International Conference on Modeling and Simulation of Microsystems*, San Juan, Puerto Rico.

- van't Veer, L., Dal, H., van de Vijver, M., He, Y., Hart, A., Mao, M., Peterse, H. van der Kooy, K., Marton, M., Witteveen, A., Schreiber, G., Kerkhoven, R., Roberts, C., Linsley, P., Bernards, R., and Friend, S. 2002. Gene expression profiling predicts clinical outcome of breast cancer. *Nature* 415, 530–536.
- Woolf, P., and Wang, Y. 2000. A fuzzy logic approach to analyzing gene expression data. *Physiol. Genomics* 3, 9–15.