

**Cached Sufficient Statistics for Efficient
Machine Learning with Large Datasets**

Andrew W. Moore and Mary Soon Lee

CMU-RI-TR-97-27

The Robotics Institute
Carnegie Mellon University
Pittsburgh, Pennsylvania 15213

July 1997

©1997 Andrew W. Moore, Mary Soon Lee

Cached Sufficient Statistics for Efficient Machine Learning with Large Datasets

Andrew W. Moore and Mary Soon Lee

Robotics Institute and
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213
(awm,mslee)@cs.cmu.edu

July 1997

Abstract

This paper introduces new algorithms and data structures for quick counting for machine learning datasets. We focus on the counting task of constructing contingency tables, but our approach is also applicable to counting the number of records in a dataset that match conjunctive queries. Subject to certain assumptions, the costs of these operations can be shown to be independent of the number of records in the dataset and loglinear in the number of non-zero entries in the contingency table.

We provide a very sparse data structure, the *ADtree*, to minimize memory use. We provide analytical worst-case bounds for this structure for several models of data distribution. We empirically demonstrate that tractably-sized data structures can be produced for large real-world datasets by (a) using a sparse tree structure that never allocates memory for counts of zero, (b) never allocating memory for counts that can be deduced from other counts, and (c) not bothering to expand the tree fully near its leaves.

We show how the *ADtree* can be used to accelerate Bayes net structure finding algorithms, rule learning algorithms, and feature selection algorithms, and we provide a number of empirical results comparing *ADtree* methods against traditional direct counting approaches. We also discuss the possible uses of *ADtrees* in other machine learning methods, and discuss the merits of *ADtrees* in comparison with alternative representations such as *kd*-trees, *R*-trees and frequent sets.

1 Caching sufficient statistics

Computational efficiency is an important concern for machine learning algorithms, especially when applied to large datasets (Fayyad et al., 1997; Fayyad and Uthurusamy, 1996) or in real-time scenarios. In (Moore et al., 1997) *kd*-trees with multiresolution cached regression

matrix statistics were used to permit very fast locally weighted and instance based regression. In this paper we attempt to accelerate predictions for symbolic attributes using a kind of *kd-tree* that splits on all dimensions at all nodes.

Many machine learning algorithms operating on datasets of symbolic attributes need to do frequent *counting*. Let us begin by establishing some notation.

We are given a data set with R records and M attributes. The attributes are called a_1, a_2, \dots, a_M . The value of attribute a_i in the k th record is a small integer lying in the range $\{1, 2, \dots, n_i\}$ where n_i is called the *arity* of attribute i . Figure 1 gives an example.

Attributes →	a ₁	a ₂	a ₃	
Arity →	n ₁ =2	n ₂ =4	n ₃ =2	
Record ₁ →	1	1	1	M = 3 R = 6
Record ₂ →	2	3	1	
Record ₃ →	2	4	2	
Record ₄ →	1	3	1	
Record ₅ →	2	3	1	
Record ₆ →	1	3	1	

Figure 1: A simple dataset used as an example. It has $R = 6$ records and $M = 3$ attributes.

Queries

A *query* is a set of (*attribute = value*) pairs in which the left hand sides of the pairs are a subset of $\{a_1 \dots a_M\}$ arranged in increasing order of index. Four examples of queries for our dataset are

$$(a_1 = 1); (a_2 = 3, a_3 = 1); (); (a_1 = 2, a_2 = 4, a_3 = 2) \quad (1)$$

Notice that the total number of possible queries is $\prod_{i=1}^M (n_i + 1)$. This is because each attribute can either appear in the query with one of the n_i values it may take, or it may be omitted (which is equivalent to giving it a $a_i = *$ “don’t care” value).

Counts

The *count* of a query, denoted by $C(\text{Query})$ is simply the number of records in the dataset matching all the (*attribute = value*) pairs in *Query*. For our example dataset we find:

$$\begin{aligned} C(a_1 = 1) &= 3 \\ C(a_2 = 3, a_3 = 1) &= 4 \\ C() &= 6 \\ C(a_1 = 2, a_2 = 4, a_3 = 2) &= 1 \end{aligned}$$

Contingency Tables

Each subset of attributes, $a_{i(1)} \dots a_{i(n)}$, has an associated *contingency table* denoted by

$ct(a_{i(1)} \dots a_{i(n)})$. This is a table with a row for each of the possible sets of values for $a_{i(1)} \dots a_{i(n)}$. The row corresponding to $a_{i(1)} = v_1 \dots a_{i(N)} = v_n$ records the count $C(a_{i(1)} = v_1 \dots a_{i(N)} = v_n)$. Our example dataset has 3 attributes and so $2^3 = 8$ contingency tables exist, depicted in Figure 2.

$ct()$		$ct(a_1)$		$ct(a_2)$		$ct(a_1, a_2, a_3)$			
#		a_1	#	a_2	#	a_1	a_2	a_3	#
6		1	3	1	1	1	1	1	1
		2	3	2	0	1	1	2	0
				3	4	1	2	1	0
				4	1	1	4	1	0
						1	4	2	0
						2	1	1	0
						2	1	2	0
						2	2	1	0
						2	2	2	0
						2	3	1	2
						2	3	2	0
						2	4	1	0
						2	4	2	1

Figure 2: The eight possible contingency tables for the dataset of Figure 1.

A *conditional contingency table*, written

$$ct(a_{i(1)} \dots a_{i(n)} \mid a_{j(1)} = u_1, \dots, a_{j(p)} = u_p) \quad (2)$$

is the contingency table for the subset of records in the dataset that match the query to the right of the \mid symbol. For example,

$$ct(a_1, a_3 \mid a_2 = 3) = \begin{array}{|c|c|c|} \hline a_1 & a_3 & \# \\ \hline 1 & 1 & 2 \\ \hline 1 & 2 & 0 \\ \hline 2 & 1 & 2 \\ \hline 2 & 2 & 0 \\ \hline \end{array}$$

Contingency tables are used in a variety of machine learning applications, including building the probability tables for Bayes nets and evaluating candidate conjunctive rules in rule learning algorithms such as (Quinlan, 1990; Clark and Niblett, 1989). It would thus be desirable to be able to perform such counting efficiently.

If we are prepared to pay a one-time cost for building a caching data structure, then it is easy to suggest a mechanism for doing counting in constant time. For each possible query, we precompute the contingency table. The total amount of numbers stored in memory for such a data structure would be $\prod_{i=1}^M (n_i + 1)$, which even for our humble dataset of Figure 1 is 45, as revealed by Figure 2. For a real dataset with more than ten attributes of medium arity, or fifteen binary attributes, this is far too large to fit in main memory.

We would like to retain the speed of precomputed contingency tables without incurring an intractable memory demand. That is the subject of this paper.

2 Cache Reduction 1: The dense ADtree for caching sufficient statistics

First we will describe the ADtree, the data structure we will use to represent the set of all possible counts. Our initial simplified description is an obvious tree representation that does not yield any immediate memory savings, but will later provide several opportunities for cutting off zero counts and redundant counts. This structure is shown in Figure 3. An ADtree node (shown as a rectangle) has child nodes called “Vary nodes” (shown as ovals). Each ADnode represents a query and stores the number of records that match the query (in the $C = \#$ field). The Vary a_j child of an ADnode has one child for each of the n_j values of attribute a_j . The k th such child represents the same query as Vary a_j ’s parent, with the additional constraint that $a_j = k$.

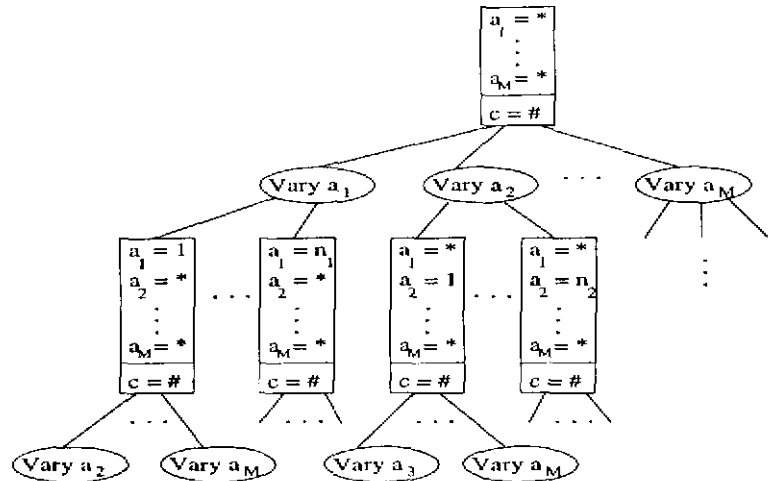


Figure 3: The top ADnodes of an ADtree, described in the text.

Notes regarding this structure:

- Although drawn on the diagram, the description of the query (e.g. $a_1 = 1, a_2 = *, \dots, a_M = *$ on the leftmost ADnode of the second level) is not explicitly recorded in the ADnode. The contents of an ADnode are simply a count and a set of pointers to the Vary a_j children.

The contents of a Vary a_j node are a set of pointers to ADnodes.

- The cost of looking up a count is proportional to the number of instantiated variables in the query. For example, to look up $C(a_7 = 2, a_{13} = 1, a_{22} = 3)$ we would follow the following path in the tree: Vary $a_7 \rightarrow a_7 = 2 \rightarrow$ Vary $a_{13} \rightarrow a_{13} = 1 \rightarrow$ Vary $a_{22} \rightarrow a_{22} = 3$. Then the count is obtained from the resulting node.
- Notice that if a node ADN has Vary a_i as its parent, then ADN’s children are

$$\text{Vary } a_{i+1} \text{ Vary } a_{i+2} \dots \text{ Vary } a_M.$$

It is not necessary to store Vary nodes with indices below $i+1$ because that information can be obtained from another path in the tree.

Cutting off nodes with counts of zero

As described, the tree is not sparse and will contain exactly $\prod_{i=1}^M (n_i + 1)$ nodes. Sparseness is easily achieved by storing a *NULL* instead of a node for any query that matches zero records. All of the specializations of such a query will also have a count of zero and they will not appear anywhere in the tree. For some datasets this can reduce the number of numbers that need to be stored. For example, the dataset in Figure 1 which previously needed 45 numbers to represent all contingency tables will now only need 22 numbers.

3 Cache Reduction II: The sparse ADtree

It is easy to devise datasets for which there is no benefit in failing to store counts of zero. Suppose we have M binary attributes and 2^M records in which the k th record is the bits of the binary representation of k . Then no query has a count of zero and the tree contains 3^M nodes.

To reduce the tree size despite this, we will take advantage of the observation that very many of the counts stored in the above tree are redundant.

Each Vary a_j node in the above ADtree stores n_j subtrees—one subtree for each value of a_j . Instead, we will find the most common of the values of a_j (call it *MCV*) and store a *NULL* in place of the *MCV*th subtree. The remaining $n_j - 1$ subtrees will be represented as before. An example for a simple dataset is given in Figure 4. Appendix 2 describes the straightforward algorithm for building such an ADtree.

As we will see in Section 4, it is still possible to build full exact contingency tables (or give counts for specific queries) in time that is only slightly larger than for the full ADtree of Section 2. But first let us examine the memory consequences of this representation.

Appendix 1 shows that for binary attributes, in the worst case, given M attributes and R records, the number of nodes needed to store the tree is bounded above by 2^M in the worst case (and much less if $R < 2^M$). In contrast, the amount of memory needed by the dense tree of Section 2 is 3^M in the worst case.

Notice in Figure 4 that the *MCV* value is context dependent. Depending on constraints on parent nodes, a_2 's *MCV* is sometimes 1 and sometimes 2. This context dependency can provide dramatic savings if (as is frequently the case) there are correlations among the attributes. This is discussed further in Appendix 1.

4 Computing contingency tables from the sparse ADtree

Given an ADtree we wish to be able to quickly construct contingency tables for any arbitrary set of attributes $\{a_{i(1)} \dots a_{i(n)}\}$.

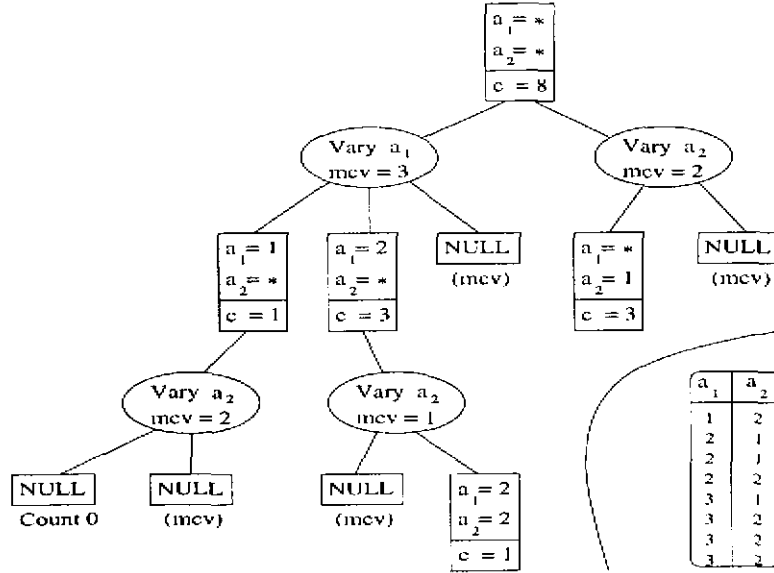


Figure 4: A sparse *ADtree* built for the dataset shown in the bottom right. The most common value for a_1 is 3, and so the $a_1 = 3$ subtree of the *Vary a_1* child of the root node is *NULL*. At each of the *Vary a_2* nodes the most common child is also set to *NULL* (which child is most common depends on the context).

Notice that a conditional contingency table $\mathbf{ct}(a_{i(1)} \dots a_{i(n)} \mid \text{Query})$ can be built recursively. We first build

$$\begin{aligned}
 & \mathbf{ct}(a_{i(2)} \dots a_{i(n)} \mid a_{i(1)} = 1, \text{Query}) \\
 & \mathbf{ct}(a_{i(2)} \dots a_{i(n)} \mid a_{i(1)} = 2, \text{Query}) \\
 & \quad \vdots \\
 & \mathbf{ct}(a_{i(2)} \dots a_{i(n)} \mid a_{i(1)} = n_{i(1)}, \text{Query})
 \end{aligned}$$

For example, to build $\mathbf{ct}(a_1, a_3)$ using the dataset in Figure 1, we can build $\mathbf{ct}(a_3 \mid a_1 = 1)$ and $\mathbf{ct}(a_3 \mid a_1 = 2)$ and combine them as in Figure 5.

When building a conditional contingency table from an *ADtree*, we will not need to explicitly specify the query condition. Instead we will supply an *ADnode* of the *ADtree*, which implicitly is equivalent information. The algorithm is:

MakeContab($\{ a_{i(1)} \dots a_{i(n)} \}$, *ADN*)
 Let $VN :=$ The *Vary $a_{i(1)}$* subnode of *ADN*.
 Let $MCV := VN.MCV$.
 For $k := 1, 2, \dots, n_{i(1)}$
 If $k \neq MCV$
 Let $ADN_k :=$ The $a_{i(1)} = k$ subnode of VN .
 $CT_k := \mathbf{MakeContab}(\{a_{i(2)} \dots a_{i(n)}\}, ADN_k)$.

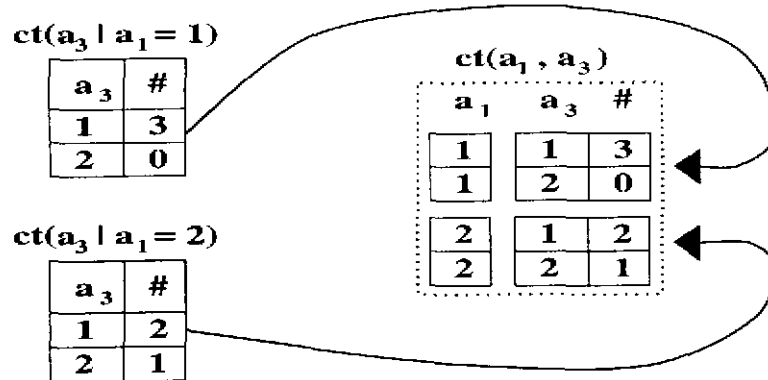


Figure 5: An example (using numbers from Figure 1) of how contingency tables can be combined recursively to form larger contingency tables.

$CT_{MCV} := ???$ (explained below)

Return the concatenation of $CT_1 \dots CT_{n_{i(1)}}$.

The base case of this recursion occurs when the first argument is empty, in which case we return a one-element contingency table containing the count associated with the current *ADnode*, *ADN*.

There is an omission in the algorithm. In the iteration over $k \in \{1, 2, \dots, n_{i(1)}\}$ we are unable to compute the conditional contingency table for CT_{MCV} because the $a_{i(1)} = MCV$ subtree is deliberately missing as per Section 3. What can we do instead?

We can take advantage of the following property of contingency tables:

$$\mathbf{ct}(a_{i(2)} \dots a_{i(n)} \mid Query) = \sum_{k=1}^{n_{i(1)}} \mathbf{ct}(a_{i(2)} \dots a_{i(n)} \mid a_{i(1)} = k, Query) \quad (3)$$

The value $\mathbf{ct}(a_{i(2)} \dots a_{i(n)} \mid Query)$ can be computed from within our algorithm by calling

$$\mathbf{MakeContab}(\{a_{i(2)} \dots a_{i(n)}\}, ADN) \quad (4)$$

and so the missing conditional contingency table in the algorithm can be computed by the following row-wise subtraction:

$$CT_{MCV} := \mathbf{MakeContab}(\{a_{i(2)} \dots a_{i(n)}\}, ADN) - \sum_{k \neq MCV} CT_k \quad (5)$$

A similar algorithm that takes Frequent Sets (see Section 8) as input and computes counts is described and evaluated in (Mannila and Toivonen, 1996).

Complexity of building a contingency table

What is the cost of computing a contingency table? Let us consider computing a contingency

table for n binary attributes. The entire contingency table has 2^n entries. Write $C(n)$ = the cost of computing such a contingency table. In the top-level call of **MakeContab** there are two calls to build contingency tables from $n - 1$ attributes: we will build $CT(a_{i(2)} \dots a_{i(n)} | a_{i(1)} = \text{LeastCommonValue}, \text{Query})$ and also $CT(a_{i(2)} \dots a_{i(n)} | \text{Query})$. Then there will be one subtraction of contingency tables, which will require 2^{n-1} numeric subtractions. So we have

$$C(0) = 1 \tag{6}$$

$$C(n) = 2C(n - 1) + 2^{n-1} \quad \text{if } n > 0 \tag{7}$$

The solution to this recurrence relation is $C(n) = (1 + n/2)2^n$ operations to build the table. By comparison, if we used no cached data structure, but simply counted through the dataset in order to build a contingency table we would need $O(nR + 2^n)$ operations where R is the number of records in the dataset. We are thus cheaper than the standard counting method if $2^{n-1} \ll R$. We are interested in large datasets in which R may be more than 100,000. In such a case our method will present a several order of magnitude speedup for, say, a contingency table of eight attributes. Notice that this cost is **independent of M**, the total number of attributes in the dataset, and only depends upon the (almost always much smaller) number of attributes n requested for the contingency table.

If the attributes all have arity k instead of being binary then the cost will be $(1 + n(k - 1)/k)k^{n-1}$. As with the binary case, this cost is also loglinear in the size of the contingency table.

Sparse representation of contingency tables

In practice we do not represent contingency tables as multidimensional arrays, but as tree structures. This gives both the slow counting approach and the *ADtree* approach a substantial computational advantage in cases where the contingency table is sparse, i.e. has many zero entries. Figure 6 shows such a sparse contingency table representation. This can mean average-case behavior is much faster than worst case for contingency tables with large numbers of attributes or high-arity attributes.

Indeed, our experiments in Section 7 show costs rising much more slowly than $O(nk^{n-1})$ as n increases. Note too that when using a sparse representation, the worst-case for **MakeContab** is now $O(\min(nR, nk^{n-1}))$ because R is the maximum possible number of non-zero contingency table entries.

5 Cache Reduction III: Leaf-lists

We now introduce a scheme for further reducing memory use. It is not worth building the *ADtree* data structure for a small number of records. For example, suppose we have 15 records and 40 binary attributes. Then the analysis in Appendix 1 shows us that in the worst case the *ADtree* might require 10701 nodes. But computing contingency tables using the resulting *ADtree* would, with so few records, be no faster than the conventional counting approach, which would merely require us to retain the dataset in memory.

Aside from concluding that *ADtrees* are not useful for very small datasets, this also leads to a final method for saving memory in large *ADtrees*. Any *ADtree* node with fewer than

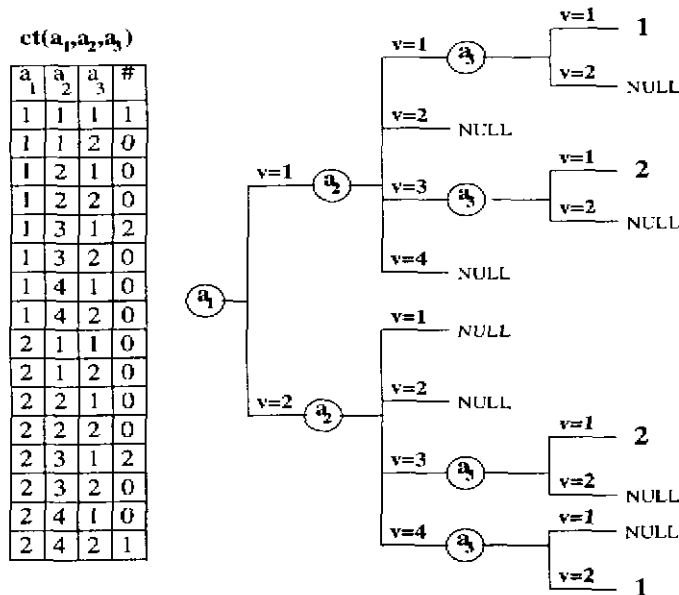


Figure 6: The right hand figure is the sparse representation of the contingency table on the left.

R_{min} records does not expand its subtree. Instead it maintains a list of pointers into the original dataset, explicitly listing those records that match the current $ADnode$. Such a list of pointers is called a *leaf-list*. Figure 7 gives an example.

The use of leaf-lists has one minor and two major consequences. The minor consequence is the need to include a straightforward change in the contingency table generating algorithm to handle leaf-list nodes. This minor alteration is not described here.

The first major consequence is that now the dataset itself must be retained in main memory so that algorithms that inspect leaf-lists can access the rows of data pointed to in those leaf-lists.

The second major consequence is that the $ADtree$ may require much less memory. This is documented in Section 7 and worst-case bounds are provided in Appendix 1.

6 Using $ADtrees$ for Machine Learning

As we will see in Section 7, the $ADtree$ structure can substantially speed up the computation of contingency tables for large real datasets. How can machine learning and statistical algorithms take advantage of this? Here we provide three examples: Feature Selection, Bayes net scoring and rule learning. But it seems likely that many other algorithms can also benefit, for example decision trees (Quinlan, 1983; Breiman et al., 1984) and GMDH (Madala and Ivakhnenko, 1994). Text classification. In future work we will also examine ways to speed up nearest neighbor and other memory-based queries using $ADtrees$.

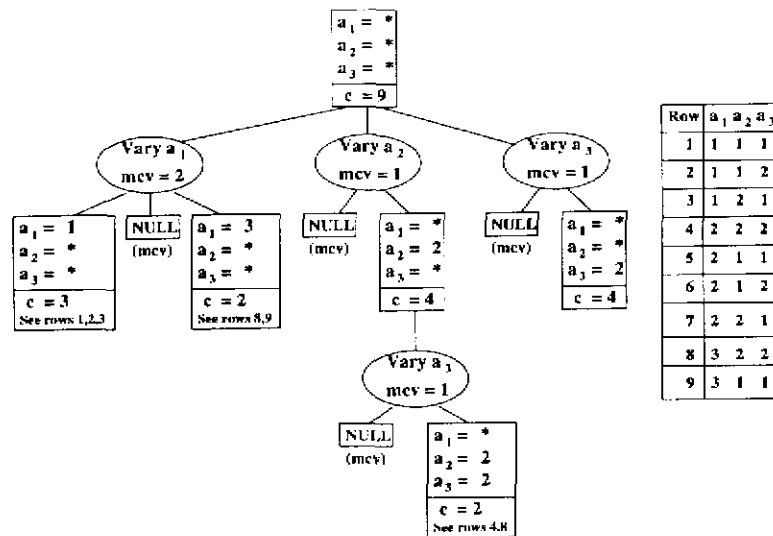


Figure 7: An *ADtree* built using leaf-lists with $R_{\min} = 4$. Any node matching 3 or fewer records is not expanded, but simply records a set of pointers into the dataset (shown on the right).

6.1 DataSets

The experiments used the following datasets.

Name	$R = \text{Num. Records}$	$M = \text{Num. Attributes}$	
ADULT1	15060	15	The small “Adult Income” dataset placed in the UCI repository by Ron Kohavi. Contains census data related to job, wealth, and nationality. Attribute arities range from 2 to 41. In the UCI repository this is called the Test Set. Rows with missing values were removed.
ADULT2	30162	15	The same kinds of records as above but with different data. The Training Set.
ADULT3	45222	15	ADULT1 and ADULT2 concatenated.
BIRTH	9672	97	Records concerning a very wide number of readings and factors recorded at various stages during pregnancy. Most attributes are binary, and 70 of the attributes are very sparse, with over 95% of the values being FALSE.
SYNTH	30K–500K	24	Synthetic datasets of entirely binary attributes generated using the Bayes net in Figure 8.

6.2 Using *ADtrees* for feature selection

Given M attributes, of which one is an output that we wish to predict, it is often interesting to ask “which subset of n attributes, ($n < M$), is the best predictor of the output on the

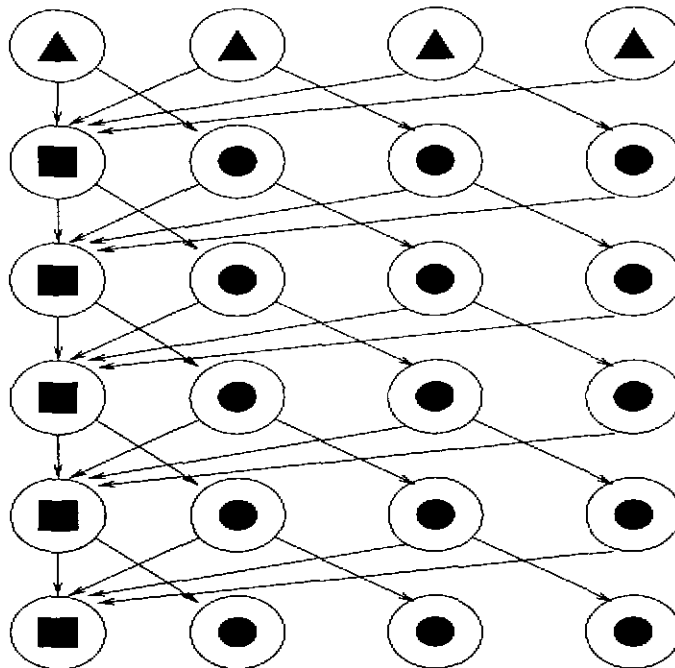


Figure 8: A Bayes net that generated our SYNTH datasets. There are three kinds of nodes. The nodes marked with triangles are generated with $P(a_i = 1) = 0.8, P(a_i = 2) = 0.2$. The square nodes are deterministic. A square node takes value 2 if the sum of its four parents is even, else it takes value 1. The circle nodes are probabilistic functions of their single parent, defined by $P(a_i = 2 | Parent = 1) = 0$ and $P(a_i = 2 | Parent = 2) = 0.4$. This provides a dataset with fairly sparse values and with many interdependencies.

same distribution of datapoints that are reflected in this dataset?” There are many ways of scoring a set of features, but a particularly simple one is *information gain* (Quinlan, 1983).

Let a_{out} be the attribute we wish to predict and let $a_{i(1)} \dots a_{i(n)}$ be the set of attributes used as inputs. Let X be the set of possible assignments of values to $a_{i(1)} \dots a_{i(n)}$ and write $Assign_k \in X$ as the k th such assignment. Then

$$InfoGain = \sum_{v=1}^{n_{out}} f\left(\frac{C(a_{out} = v)}{R}\right) - \sum_{k=1}^{|X|} \frac{C(Assign_k)}{R} \sum_{v=1}^{n_{out}} f\left(\frac{C(a_{out} = v, Assign_k)}{C(Assign_k)}\right) \quad (8)$$

where R is the number of records in the entire dataset and

$$f(x) = -x \log_2 x \quad (9)$$

The counts needed in the above computation can be read directly from $ct(a_{out}, a_{i(1)} \dots a_{i(n)})$.

Searching for the best subset of attributes is simply a question of search among all attribute-sets of size n (n specified by the user). This is a simple example designed to test our counting methods: any practical feature selector would need to penalize the number of rows in the contingency table (else high arity attributes would tend to win).

education-num	sex	<=50K	>50K
below6	Female	617	10
below6	Male	1562	127
6to12	Female	6434	522
6to12	Male	10183	3119
above12	Female	1619	580
above12	Male	2239	3150

Entropy(class) = 0.809566
IG(class | education-num sex) = 0.114711

Figure 9: Output from the feature selector

Figure 9 gives an example of the output of this feature selector when applied to the task of finding the two best features with arity less than 6 in the ADULT2 dataset for predicting the `class` attribute (Kohavi, 1996).

6.3 Using ADtrees for Bayes net structure discovery

There are many possible Bayes net learning tasks, all of which entail counting, and hence might be speeded up by ADtrees. In this paper we present experimental results for the particular example of scoring the structure of a Bayes net to decide how well it matches the data.

We will use maximum likelihood scoring with a penalty for the number of parameters. We first compute the probability table associated with each node. Write $Parents(j)$ for the parent attributes of node j and write X_j as the set of possible assignments of values to $Parents(j)$. The maximum likelihood estimate for

$$P(a_j = v \mid X_j) \tag{10}$$

is estimated as

$$\frac{C(a_j = v, X_j)}{C(X_j)} \tag{11}$$

and all such estimates for node j 's probability tables can be read from $ct(a_j, Parents(j))$.

The next step in scoring a structure is to decide the likelihood of the data given the probability tables we computed and to penalize the number of parameters in our network (without the penalty the likelihood would increase every time a link was added to the network). Following (Friedman and Yakhini, 1996) the penalized log-likelihood score is

$$- N_{\text{params}} \log(R)/2 + R \sum_{j=1}^M \sum_{Asgn \in X_j} \sum_{v=1}^{n_j} P(a_j = v \wedge Asgn) \log P(a_j = v \mid Asgn) \tag{12}$$

where N_{params} is the total number of probability table entries in the network.

attribute	score	np	
age	0.908279	2	pars = <no parents>
workclass	1.39105	18	pars = age
fnlwgt	0.0410872	2	pars = <no parents>
education	2.87019	105	pars = workclass
education-num	0.0	32	pars = education
marital-status	1.68602	54	pars = age education-num
occupation	2.9862	273	pars = workclass education-num
relationship	1.06944	105	pars = education-num marital-status
race	0.756854	28	pars = marital-status
sex	0.430973	84	pars = occupation relationship
capital-gain	0.041708	6	pars = education-num
capital-loss	0.233887	12	pars = relationship
hours-per-week	0.855615	84	pars = workclass education-num sex
native-country	0.69222	200	pars = race
class	0.559691	54	pars = education-num relationship
hours-per-week			
Score is 445928			
The search took 252 seconds.			

Figure 10: Output from the Bayes structure finder running on the ADULT2 dataset. *Score* is the contribution to the sum in Equation 12 due to the specified attribute. *np* is the number of entries in the probability table for the specified attribute.

We search among structures to find the best score. In these experiments we use random-restart stochastic hillclimbing in which the operations are random addition or removal of a network link. Only the probability table of the affected node is recomputed on each step. The search is restricted to look only for networks in which the indices of node j 's parents must all be less than j .

Figure 10 shows the Bayes net structure returned by our Bayes net structure finder after 100,000 iterations of hillclimbing. When running on the synthetic datasets a nearly perfect structure was typically found after a few hundred thousand iterations.

6.4 Using ADtrees for rule finding

Given an output attribute a_{out} and a distinguished value v_{out} , rule finders search among conjunctive queries of the form

$$Assign = (a_{i(1)} = v_1 \dots a_{i(N)} = v_n) \quad (13)$$

to find the query that maximizes the estimated value

$$P(a_{out} = v_{out} \mid Assign) = \frac{C(a_{out} = v_{out}, Assign)}{C(Assign)} \quad (14)$$

```

Rule 0: score = 0.965 (218/226), workclass = Private, education-num = above12,
marital-status = Married-civ-spouse, capital-loss = above1600
Rule 1: score = 0.961 (199/207), workclass = Private, education-num = above12,
relationship = Husband, capital-loss = above1600
Rule 2: score = 0.955 (193/202), age = below45, education-num = above12,
relationship = Husband, capital-loss = above1600
Rule 3: score = 0.950 (302/318), education-num = above12, relationship = Husband,
capital-loss = above1600, native-country = United-States
Rule 4: score = 0.948 (201/212), education-num = above12, relationship = Husband,
capital-loss = above1600, hours-per-week = above40

```

Figure 11: Output from the rule finder

To avoid rules without significant support, we also insist that $C(\text{Assign})$ (the number of records matching the query) must be above some threshold S_{\min} .

In these experiments we implemented a brute force search that looks through all possible queries that involve a user-specified number of attributes, n . We build each $\text{ct}(a_{out}, a_{i(1)} \dots a_{i(n)})$ in turn (there are $\binom{M}{n}$ such tables), and then look through the rows of each table for all queries using the $a_{i(1)} \dots a_{i(n)}$ that have greater than minimum support S_{\min} . We return a priority queue of the highest scoring rules.

Figure 11 gives the five highest-scoring rules found for the ADULT2 dataset using rules of size 4.

7 Experimental Results

Let us first examine the memory required by an *ADtree* on our datasets. Table 1 shows us, for example, that the ADULT2 dataset produced an *ADtree* with 95,000 nodes. The tree required almost 11 Mbytes of memory. Among the three ADULT datasets, the size of the tree varied approximately linearly with the number of records.

Unless otherwise specified, in all the experiments in this section, the ADULT datasets used no leaf-lists. The BIRTH and SYNTHETIC datasets used leaf-lists of size $R_{\min} = 16$ by default. The BIRTH dataset with its large number of sparse attributes required a modest 8 Mbytes to store the tree—many magnitudes below the worst-case bounds. Among the synthetic datasets, the tree size increased sublinearly with the dataset size. This indicates that as the dataset gets larger, novel records (which may cause new nodes to appear in the tree) become less frequent.

Table 2 shows the costs of performing 100,000 iterations of Bayes net structure searching. All experiments were performed on a 200Mhz Pentium Pro machine with 192 Mbytes of main memory. Recall that each Bayes net iteration involves one random change to the network and so requires recomputation of one contingency table. (The exception is the first iteration, in which all nodes must be computed). This means that the time to run 100,000 iterations is essentially the time to compute 100,000 contingency tables. Among the ADULT datasets, the advantage of the *ADtree* over conventional counting ranges between a factor of 16 to 35. Unsurprisingly, the computational costs for ADULT increase sublinearly with dataset size

Dataset	M	R	Nodes	Megabytes	Build Time
ADULT1	15	15060	58200	7.0	6
ADULT2	15	30162	94900	10.9	10
ADULT3	15	45222	162900	15.5	15
BIRTH	97	9672	87400	7.9	14
SYN30K	24	30000	34100	2.8	8
SYN60K	24	60000	59300	4.9	17
SYN125K	24	125000	95400	7.9	36
SYN250K	24	250000	150000	12.4	73
SYN500K	24	500000	219000	18.2	150

Table 1: The size of *ADtrees* for various datasets. *M* is the number of attributes. *R* is the number of records. *Nodes* is the number of nodes in the *ADtree*. *Megabytes* is the amount of memory needed to store the tree. *Build Time* is the number of seconds needed to build the tree (to the nearest second).

Dataset	M	R	<i>ADtree</i> Time	Regular Time	Speedup Factor
ADULT1	15	15060	228	3844	16.9
ADULT2	15	30162	316	7860	24.9
ADULT3	15	45222	348	11968	34.4
BIRTH	97	9672	44	3480	79.1
SYN30K	24	30000	24	10088	420.3
SYN60K	24	60000	24	19880	828.3
SYN125K	24	125000	24	41552	1731.3
SYN250K	24	250000	24	84496	3520.7
SYN500K	24	500000	24	168544	7022.7

Table 2: The time (in seconds) to perform 100,000 hill-climbing iterations searching for the best Bayes net structure. *ADtree Time* is the time when using the *ADtree* and *Regular Time* is the time taken when using the conventional probability table scoring method of counting through the dataset. *Speedup Factor* is the number of times by which the *ADtree* method is faster than the conventional method.

for the *ADtree* but linearly for the conventional counting. The computational advantages and the sublinear behavior are much more pronounced for the synthetic data.

Next, Table 3 examines the effect of leaf-lists on the ADULT2 and BIRTH datasets. For the ADULT dataset, the byte size of the tree decreases by a factor of 5 when leaf-lists are increased from 1 to 64. But the computational cost of running the Bayes search increases by only 25%, indicating a worth-while tradeoff if memory is scarce.

The Bayes net scoring results involved the average cost of computing contingency tables of many different sizes. The following results in Tables 4 and 5 make the savings for fixed size attribute sets easier to discern. These tables give results for the feature selection and rule finding algorithms respectively. The biggest savings come from small attribute sets. Computational savings for sets of size one or two are, however, not particularly interesting since all such counts could be cached by straightforward methods without needing any tricks. In all cases, however, we do see large savings, especially for the BIRTH data. Datasets with

R_{\min}	ADULT2				BIRTH			
	#Mb	#nodes	Build Secs	Search Secs	#Mb	#nodes	Build Secs	Search Secs
1	10.87	94,872	13	316				
2	8.46	86,680	10	312				
4	6.30	75,011	8	308	27.60	245,722	31	40
8	4.62	62,095	7	312	14.90	152,409	22	40
16	3.37	49,000	7	320	7.95	87,415	15	44
32	2.55	37,790	5	344	4.30	46,777	11	44
64	1.98	27,726	4	388	2.42	23,474	8	52
128	1.59	18,903	4	472	1.48	11,554	6	64
256	1.38	12,539	4	600	0.95	5,150	5	92
512	1.21	7,336	2	864	0.65	2,237	3	148
1024	1.05	3,938	2	1344	0.45	887	3	244
2048	0.95	1,890	2	1856	0.28	246	2	364
4096	0.86	780	2	2808	0.28	201	1	380

Table 3: Investigating the effect of the R_{\min} parameter on the ADULT2 dataset and the BIRTH dataset. *#Mb* is the memory used by the *ADtree*. *#nodes* is the number of nodes in the *ADtree*. *Build Secs* is the time to build the *ADtree*. *Search Secs* is the time needed to perform 100,000 iterations of the Bayes net structure search.

larger numbers of rows would, of course, reveal larger savings.

8 Alternative Data Structures

Why not use a *kd*-tree?

kd-trees can be used for accelerating learning algorithms (Omohundro, 1987; Moore et al., 1997). The primary difference is that a *kd*-tree node splits on only one attribute instead of all attributes. This results in much less memory (linear in the number of records). But

Number Attributes	ADULT2				BIRTH			
	Number Attribute Sets	<i>ADtree</i> Time	Regular Time	Speedup Factor	Number Attribute Sets	<i>ADtree</i> Time	Regular Time	Speedup Factor
1	14	.000071	.048	675.0	96	.000018	.015	841
2	91	.00054	.067	124.0	4,560	.000042	.021	509
3	364	.0025	.088	34.9	142,880	.000093	.028	298
4	1,001	.0083	.11	13.4	3,321,960	.00019		
5	2,002	.023	.14	6.0	61,124,064	.00033		

Table 4: The time taken to search among all attribute sets of a given size (Number Attributes) for the set that gives the best information gain in predicting the output attribute. The times, in seconds, are the average evaluation times per attribute-set.

Number Attributes	ADULT2				BIRTH			
	Number Rules	<i>ADtree</i> Time	Regular Time	Speedup Factor	Number Rules	<i>ADtree</i> Time	Regular Time	Speedup Factor
1	116	.000019	.0056	295.0	194	.000025	.0072	286
2	4,251	.000019	.0014	75.3	17,738	.000021	.0055	259
3	56,775	.000024	.00058	23.8	987,134	.000022	.0040	186
4	378,984	.000031	.00030	9.8	37,824,734	.000024	.0030	127
5	1,505,763	.000042	.00019	4.7	1,077,672,055	.000026		

Table 5: The time taken to search among all rules of a given size (Number Attributes) for the highest scoring rules for predicting the output attribute. The times, in seconds, are the average evaluation time per rule.

counting can be expensive. Suppose, for example, that level one of the tree splits on a_1 , level two splits on a_2 etc. Then in the case of binary variables, if we have a query involving only attributes a_{20} and higher then we have to explore all paths in the tree down to level 20. With datasets of fewer than 2^{20} records this may be no cheaper than performing a linear search through the records. Another possibility are R -trees (Guttman, 1984; Roussopoulos and Leifker, 1985), which store databases of M -dimensional geometric objects. However, in this context, they offer no advantages over kd -trees.

Why not use a frequent set finder?

Frequent set finders (Agrawal et al., 1996) are typically used with very large databases of millions of records containing very sparse binary attributes. Efficient algorithms exist for finding all subsets of attributes that co-occur with value TRUE in more than a fixed number (chosen by the user, and called the *support*) of records. (Mannila and Toivonen, 1996) suggests that such frequent sets can be used to perform efficient counting. In the case where *support* = 1, all such frequent sets are gathered and if counts of each frequent set are retained, this is equivalent to producing an *ADtree* in which instead of performing a node cutoff for the most common value, the cutoff always occurs for value FALSE.

The use of frequent sets in this way would thus be very similar to the use of *ADtrees*, with one advantage and one disadvantage. The advantage is that efficient algorithms have been developed for building frequent sets from a small number of sequential passes through data. The *ADtree* requires random access to the dataset while it is being built, and for its leaf-lists. This is impractical if the dataset is too large to reside in main memory and is accessed through database queries.

The disadvantage of frequent sets in comparison with *ADtrees* is that under some circumstances they may require much more memory. Assume the value 2 is rarer than 1 throughout all attributes in the dataset and assume reasonably that we thus choose to find all frequent sets of 2's. Unnecessarily many sets will be produced if there are correlations. In the extreme case, imagine a dataset in which 30% of the values are 2, 70% are 1 and attributes are perfectly correlated—all values in each record are identical! Then with M attributes there would be 2^M frequent sets of 2's. In contrast, the *ADtree* would only contain $M + 1$ nodes. This is an extreme example, but datasets with much weaker inter-attribute correlations can

similarly benefit from using an *ADtree*.

Leaf-lists are another technique to reduce the size of *ADtrees* further. They could also be used for the frequent set representation.

9 Discussion

What about numeric attributes?

The *ADtree* representation is designed entirely for symbolic attributes. When faced with numeric attributes, the simplest solution is to discretize them into a fixed finite set of values which are then treated as symbols, but *this is of little help if the user requests counts for queries involving inequalities on numeric attributes*. In future work we will evaluate the use of structures combining elements from multiresolution *kd-trees* of real attributes (Moore et al., 1997) with *ADtrees*.

Algorithm-specific counting tricks

Many algorithms that count using the conventional “linear” method have algorithm-specific ways of accelerating their performance. For example, a Bayes net structure finder may try to remember all the contingency tables it has tried previously in case it needs to re-evaluate them. When it deletes a link, it can deduce the new contingency table from the old one without needing a linear count.

In such cases, the most appropriate use of the *ADtree* may be as a lazy caching mechanism. At birth, the *ADtree* consists only of the root node. Whenever the structure finder needs a contingency table that cannot be deduced from the current *ADtree* structure, the appropriate nodes of the *ADtree* are expanded. The *ADtree* then takes on the role of the algorithm-specific caching methods, while (in general) using up much less memory than if all contingency tables were remembered.

Hard to update incrementally

Although the tree can be built cheaply (see the experimental results in Section 7), and although it can be built lazily, the *ADtree* cannot be updated cheaply with a new record. This is because one new record may match up to 2^M nodes in the tree in the worst case.

Scaling up

The *ADtree* representation can be useful for datasets of the rough size and shape used in this paper. But they cannot represent all the sufficient statistics for huge datasets with many hundreds of non-sparse and poorly correlated attributes. What should we do if our dataset or our *ADtree* cannot fit into main memory? In the latter case, we could simply increase the size of leaf-lists, trading off decreased memory against increased time to build contingency tables. But if that is inadequate at least three possibilities remain. First, we could build approximate *ADtrees* that don’t store any information for nodes that match fewer than a threshold number of records. Then approximate contingency tables (complete with error bounds) can be produced using methods described in (Mannila and Toivonen, 1996). A second possibility is to exploit secondary storage and store deep, rarely visited nodes of the

ADtree on disk. This would doubtless best be achieved by integrating the machine learning algorithms with current database management tools—a topic of considerable interest in the data mining community (Fayyad et al., 1997). A third possibility, which restricts the size of contingency tables we may ask for, is to refuse to store counts for queries with more than some threshold number of attributes.

***ADtrees* are algorithm and learning-task independent**

Once an *ADtree* is built, it may be saved to disk along with the dataset and reused whenever an algorithm needs contingency tables. The tree does not treat input and output attributes any differently, so the same tree may be used for multiple tasks. Finally, a learning algorithm might wish to run over a restricted subset of the data specified by a query. For example we might want to run a rule learner restricted to the subset of records in which $a_1 = 1$ and $a_7 = 3$. In that case the same *ADtree* can also be re-used.

Acknowledgements

This work was sponsored by a National Science Foundation Career Award to Andrew Moore. The authors thank Justin Boyan, Scott Davies and Jeff Schneider for their suggestions.

Appendix 1: Memory Costs

In this appendix we examine the size of the tree. For simplicity, we restrict attention to the case of binary attributes.

The worst-case number of nodes in an *ADtree*

Given a dataset with M attributes and R records, the worst-case for the *ADtree* will occur if all 2^M possible records exist in the dataset. Then, for every subset of attributes there exists exactly one node in the *ADtree*. For example consider the attribute set $\{a_{i(1)} \dots a_{i(n)}\}$, where $i(1) < i(2) < \dots < i(n)$. Suppose there is a node in the tree corresponding to the query $\{a_{i(1)} = v_1 \dots a_{i(n)} = v_n\}$ for some values $v_1 \dots v_n$. From the definition of an *ADtree*, and remembering we are only considering the case of binary attributes, we can state:

- v_1 is the least common value of $a_{i(1)}$.
- v_2 is the least common value of $a_{i(2)}$ among those records that match $(a_{i(1)} = v_1)$.
- \vdots
- v_{k+1} is the least common value of $a_{i(k+1)}$ among those records that match $(a_{i(1)} = v_1, \dots, a_{i(k)} = v_k)$.

So there is at most one such node. Moreover since our worst-case assumption is that all possible records exist in the database, we see that the *ADtree* will indeed contain this node. Thus, the worst-case number of nodes is the same as the number of possible subsets of attributes: 2^M .

The worst-case number of nodes in an *ADtree* with a reasonable number of rows

It is frequently the case that a dataset has $R \ll 2^M$. With fewer records, there is a much lower worst-case bound on the *ADtree* size. A node at the k th level of the tree corresponds to a query involving k attributes (counting the root node as level 0). Such a node can match at most $R2^{-k}$ records because each of the node's ancestors up the tree has pruned off at least half the records by choosing to expand only the least common value of the attribute introduced by that ancestor. Thus there can be no tree nodes at level $\lfloor \log_2 R \rfloor + 1$ of the tree, because such nodes would have to match fewer than $R2^{-\lfloor \log_2 R \rfloor - 1} < 1$ records. They would thus match no records, making them *NULL*.

The nodes in an *ADtree* must all exist at level $\lfloor \log_2 R \rfloor$ or higher. The number of nodes at level k is at most $\binom{M}{k}$, because every node at level k involves an attribute set of size k and because (given binary attributes) for every attribute set there is at most one node in the *ADtree*. Thus the total number of nodes in the tree, summing over the levels is less than

$$\sum_{k=0}^{\lfloor \log_2 R \rfloor} \binom{M}{k} \text{ bounded above by } O(M^{\lfloor \log_2 R \rfloor} / (\lfloor \log_2 R \rfloor - 1)!) \quad (15)$$

The number of nodes assuming skewed independent attribute values

Imagine that all values of all attributes in the dataset are independent random binary variables, taking value 2 with probability p and taking value 1 with probability $1 - p$. Then the further p is from 0.5, the smaller we can expect the *ADtree* to be. This is because, on average, the less common value of a Vary node will match fraction $\min(p, 1 - p)$ of its parent's records. And, on average, the number of records matched at the k th level of the tree will be $R(\min(p, 1 - p))^k$. Thus, the maximum level in the tree at which we may find a node matching one or more records is approximately $\lfloor (\log_2 R) / (-\log_2 q) \rfloor$, where $q = \min(p, 1 - p)$. And so the total number of nodes in the tree is approximately

$$\sum_{k=0}^{\lfloor (\log_2 R) / (-\log_2 q) \rfloor} \binom{M}{k} \text{ bounded above by } O(M^{\lfloor (\log_2 R) / (-\log_2 q) \rfloor} / (\lfloor (\log_2 R) / (-\log_2 q) \rfloor - 1)!) \quad (16)$$

Since the exponent is reduced by a factor of $\log_2(1/q)$, skewedness among the attributes thus brings enormous savings in memory.

The number of nodes assuming correlated attribute values

The *ADtree* benefits from correlations among attributes in much the same way that it benefits from skewedness. For example, suppose that each record was generated by the simple Bayes net in Figure 12, where the random variable B is hidden (not included in the record). Then for $i \neq j$, $P(a_i \neq a_j) = 2p(1 - p)$. If *ADN* is any node in the resulting *ADtree* then the number of records matching any other node two levels below *ADN* in the tree will be fraction $2p(1 - p)$ of the number of records matching *ADN*. From this we can see that the number of nodes in the tree is approximately

$$\sum_{k=0}^{\lfloor (\log_2 R) / (-\log_2 q) \rfloor} \binom{M}{k} \text{ bounded above by } O(M^{\lfloor (\log_2 R) / (-\log_2 q) \rfloor} / (\lfloor (\log_2 R) / (-\log_2 q) \rfloor - 1)!) \quad (17)$$

where $q = \sqrt{2p(1-p)}$. Correlation among the attributes can thus also bring enormous savings in memory even if (as is the case in our example) the marginal distribution of individual attributes is uniform.

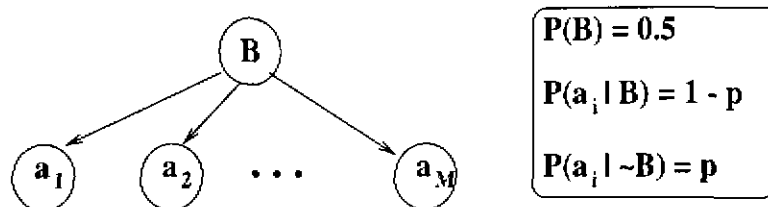


Figure 12: A Bayes net that generates correlated boolean attributes $a_1, a_2 \dots a_M$.

The number of nodes for the dense ADtree of Section 2

The dense ADtrees do not cut off the tree for the most common value of a Vary node. The worst case ADtree will occur if all 2^M possible records exist in the dataset. Then the dense ADtree will require 3^M nodes because every possible query (with each attribute taking values 1, 2 or *) will have a count in the tree. The number of nodes at the k th level of the dense ADtree can be $2^k \binom{M}{k}$ in the worst case.

The number of nodes when using Leaf-lists

Leaf-lists were described in Section 7. If a tree is built using maximum leaf-list size of R_{\min} , then any node in the ADtree matching fewer than R_{\min} records is a leaf node. This means that formulae 15, 16 and 17 can be re-used, replacing R with R/R_{\min} . It is important to remember, however, that the leaf nodes must now contain room for R_{\min} numbers instead of a single count.

Appendix 2: Building the ADtree

We define the function **MakeADTree**(a_i , $RecordNums$) where $RecordNums$ is a subset of $\{1, 2, \dots, R\}$ (R is the total number of records in the dataset) and where $1 \leq i \leq M$. This makes an ADtree from the rows specified in $RecordNums$ in which all ADnodes represent queries in which only attributes a_i and higher are used.

```

MakeADTree( $a_i$ ,  $RecordNums$ )
  Make a new ADnode called  $ADN$ .
   $ADN.COUNT := |RecordNums|$ .
  For  $j := i, i + 1, \dots, M$ 
     $j$ th Vary node of  $ADN := \mathbf{MakeVaryNode}(a_j, RecordNums)$ .

```

```

MakeVaryNode( $a_i$ ,  $RecordNums$ )
  Make a new Vary node called  $VN$ .
  For  $k := 1, 2, \dots, n_i$ 
    Let  $Childnums_k := \{\}$ .
  For each  $j \in RecordNums$ 
    Let  $v_{ij} =$  Value of attribute  $a_i$  in record  $j$ 
    Add  $j$  to the set  $Childnums_{v_{ij}}$ 
  Let  $VN.MCV := \operatorname{argmax}_k |Childnums_k|$ .
  For  $k := 1, 2, \dots, n_i$ 
    If  $|Childnums_k| = 0$  or if  $k = MCV$ 
      Set the  $a_i = k$  subtree of  $VN$  to NULL.
    Else
      Set the  $a_i = k$  subtree of  $VN$  to MakeADTree( $a_{i+1}, Childnums_k$ )

```

To build the entire tree, we must call **MakeADTree**($a_1, \{1 \dots R\}$). Assuming binary attributes, the cost of building a tree from R records and M attributes is bounded above by

$$\sum_{k=0}^{\lceil \log_2 R \rceil} \frac{R}{2^k} \binom{M}{k} \quad (18)$$

References

- Agrawal, R., Mannila, H., Srikant, R., Toivonen, H., and Verkamo, A. I. (1996). Fast discovery of association rules. In Fayyad, U. M., Piatetsky-Shapiro, G., Smyth, P., and Uthurusamy, R., editors, *Advances in Knowledge Discovery and Data Mining*. AAAI Press.
- Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984). *Classification and Regression Trees*. Wadsworth.
- Clark, P. and Niblett, R. (1989). The CN2 induction algorithm. *Machine Learning*, 3:261–284.
- Fayyad, U., Mannila, H., and Piatetsky-Shapiro, G. (1997). *Data Mining and Knowledge Discovery*. Kluwer Academic Publishers. A new journal.
- Fayyad, U. and Uthurusamy, R. (1996). Special issue on Data Mining. *Communications of the ACM*, 39(11).
- Friedman, N. and Yakhini, Z. (1996). On the sample complexity of learning Bayesian networks. In *Proceedings of the 12th conference on Uncertainty in Artificial Intelligence*.

- Guttman, A. (1984). R-trees: A dynamic index structure for spatial searching. In *Proceedings of SIGMOD 84*.
- Kohavi, R. (1996). Scaling up the accuracy of naive-Bayes classifiers: a decision-tree hybrid. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*.
- Madala, H. R. and Ivakhnenko, A. G. (1994). *Inductive Learning Algorithms for Complex Systems Modeling*. CRC Press Inc., Boca Raton.
- Mannila, H. and Toivonen, H. (1996). Multiple uses of frequent sets and condensed representations. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*.
- Moore, A. W., Schneider, J., and Deng, K. (1997). Efficient Locally Weighted Polynomial Regression Predictions. In *Proceedings of the 1997 International Machine Learning Conference*. Morgan Kaufmann.
- Omohundro, S. M. (1987). Efficient Algorithms with Neural Network Behaviour. *Journal of Complex Systems*, 1(2):273–347.
- Quinlan, J. R. (1983). Learning Efficient Classification Procedures and their Application to Chess End Games. In Michalski, R. S., Carbonell, J. G., and Mitchell, T. M., editors, *Machine Learning—An Artificial Intelligence Approach (I)*. Tioga Publishing Company, Palo Alto.
- Quinlan, J. R. (1990). Learning logical definitions from relations. *Machine Learning*, 5:239–266.
- Roussopoulos, N. and Leifker, D. (1985). Direct spatial search on pictorial databases using packed R-trees. In *ACM SIGMOD 1985 Proceedings*.