

Case-Factor Diagrams for Structured Probabilistic Modeling

David McAllester

TTI at Chicago
mcallester@tti-c.org

Michael Collins

CSAIL
Massachusetts Institute of Technology
mcollins@ai.mit.edu

Fernando Pereira

CIS
University of Pennsylvania
pereira@cis.upenn.edu

April 11, 2005

Abstract

We introduce a probabilistic formalism subsuming Markov random fields of bounded tree width and probabilistic context free grammars. Our models are based on a representation of Boolean formulas that we call case-factor diagrams (CFDs). CFDs are similar to binary decision diagrams (BDDs) but are more concise than BDDs for circuits of bounded tree width and can concisely represent the set of parse trees over a given string under a given context free grammar (unlike BDDs). A probabilistic model consists of a CFD defining a feasible set of Boolean assignments and a weight (or cost) for each individual Boolean variable. We give an inside-outside algorithm for simultaneously computing the marginal of each Boolean variable, and a Viterbi algorithm for finding the minimum cost variable assignment. Both algorithms run in time proportional to the size of the CFD.

1 Introduction

In this paper, we investigate efficient representations for structured probabilistic models. Informally, a structured model defines a distribution on structured objects such as sequences, parse trees, or assignments of values to variables. The number of possible structured objects typically grows exponentially in a natural measure of problem size. For example, the number of possible parse trees grows exponentially in the length of the string being parsed. Structured statistical models include Markov random fields (MRFs), probabilistic context free grammars (PCFGs), hidden Markov models (HMMs), conditional random fields (CRFs) [17], dynamic Bayes nets [16], probabilistic Horn abduction [20], and probabilistic relational models (PRMs) [14].

For each of these model types one can define a corresponding structured classification problem. In HMMs, for example, the problem is to recover the hidden state sequence from the observable sequence. For PCFGs, the problem is to recover a parse tree from a given word string. In PRMs, the problem is to recover latent entity labels and relations for a given set of observed entities and relations. We follow an approach where the statistical model defines $P(y|x)$ and structured classification finds the most likely y for a given x . (Other approaches are possible – for example, maximum margin classifiers are discussed below.)

The structured statistical models discussed above are intuitively similar. They all involve local probability tables or local cost functions. It is widely believed that many, if not all, of the above modeling formalisms can be viewed as special cases of MRFs (undirected graphical models). More specifically, in a structured classification problem one should be able to represent $P(y|x)$ as an MRF. By assuming $P(y|x)$ is modeled as an MRF one can prove theorems and design algorithms and software at an abstract level which simultaneously applies to all of the modeling formalisms discussed above.

Unfortunately, for some of the above models the representation of $P(y|x)$ as an MRF is problematic. The most problematic case is perhaps PCFGs. It is fairly easy to construct an MRF representing $P(y|x)$ where y is a parse tree and x is a word string. Unfortunately, standard MRF algorithms take exponential time when applied to the natural MRF representation. This is a somewhat surprising outcome, given that there are well-known inference algorithms for PCFGs which run in cubic time in the length of the word string x .

This paper presents a modeling formalism which handles both MRFs of bounded tree width and PCFGs. First we define a *linear Boolean model* (LBM). An LBM consists of three parts: a set of boolean variables; a formula defining a set of possible assignments to these variables (a “feasible set”); and an assignment of a weight to each variable. The weight for a complete variable assignment is then the sum of weights for those variables in the assignment that are true. The weight associated with a truth assignment can be written as a linear func-

tion of the bits in the assignment — hence the term “linear”. We show how to encode both standard MRFs and PCFGs, as LBMs.

The main problem we solve is how to encode compactly the set of possible assignments to the variables in an LBM in a single formalism handling both MRFs of bounded tree width and PCFGs. The *case-factor diagrams* (CFDs) we introduce for that purpose are similar to binary-decision diagrams (BDDs) [6]. CFDs differ from BDDs in two ways. First, CFDs are similar to *zero-suppressed* BDDs (ZBDDs) [19]. ZBDDs are designed for sparse truth assignments — truth assignments where most of the Boolean variables are false. Sparseness is important for representing PCFGs. In addition to being zero-suppressed, CFDs have “factor nodes” which allow a concise representation of problems that factor into independent subproblems. Factoring is important for representing MRFs of bounded tree width. We describe algorithms for CFDs that compute partition functions under Gibbs distributions for $P(y|x)$, that select the maximum likelihood (Viterbi) structure, and an inside-outside algorithm for computing the marginal distributions of all of the Boolean variables. These algorithms all run in time linear in the number of nodes in the CFD. We demonstrate that PCFG models can be encoded in a CFD which has $O(n^3)$ size where n is the length of the input string. We also show that MRFs of bounded tree width can be represented by a CFD with a polynomial number of nodes.

There are various lines of related work. A variant of BDDs for circuits of bounded tree width was introduced by McMillan [18]. Although McMillan’s formalism is more elaborate, it turns out that simply extending BDDs with “and” nodes suffices for representing MRFs of bounded width. But representing PCFGs seems to require a zero-suppressed formalism. Case-factor diagrams combine zero suppression with factoring — a combination that seems essential to efficient representation of PCFG parsing problems.

CFDs are closely related to the recursive conditioning algorithm introduced by Darwiche [8, 1]. Recursive conditioning cases on the value of a variable; factors the remaining problem into independent subproblems; and then solves the subproblems recursively. The nodes of a CFD correspond to the “subproblems” that arise in recursive conditioning. Darwiche has also defined a data structure for representing the subproblems of recursive conditioning based on arithmetic expressions [9]. The differences between CFDs and Darwiche’s expressions involve the generalization to the problem of parsing PCFGs — a problem not addressed by Darwiche. CFDs and recursive conditioning can both exploit context-sensitive independence [5] where two variables are independent under some values of a third variable but not independent in general. Context sensitive independence is particularly important for PCFGs where the tree width of the natural MRF representation is large. However, as is explained in section 8, getting $O(n^3)$ time behavior rather than $O(n^6)$ time behavior seems to require handling context sensitive *variable existence* as well as context sensitive independence. In parsing, although there are $\Omega(n^2)$ possible substrings of the

given input string, only $O(n)$ of those substrings represent phrases in any single parse and a parse is determined by the value of only $O(n)$ choices. This is very much unlike a Bayesian network where all variables are used in all assignments. In CFDs context sensitive variable existence is handled with zero suppression. Boolean variables that are forced to be false need not be mentioned — in certain contexts these variables essentially cease to exist.

Rina Dechter has given an and-or graph data structure similar to CFDs and to Darwiche’s expression graphs for representing Bayesian networks [10]. Another closely related formalism has been given by Jaeger [15] and related algorithms similar to recursive conditioning have been given by Bacchus, Dalmo and Pitassi [3, 4]. As with Darwiche’s formalism, however, these formalisms do not address the need for context sensitive variable existence in the PCFG parsing problem.

It is also important to note that CFDs are similar to BDDs in that they carry a semantics independent of the way in which they were constructed. An independent semantics allows one to treat CFDs as first class models. One of the fundamental properties of BDDs is that one can define Boolean and projection operations on BDDs allowing them to be built up in a compositional way. CFDs, as a form of BDD, can also support compositional operations although we do not explore such operations here.

Developing a common language for structured modeling has potential applications to maximum-margin structured classification [22, 7, 2]. A maximum margin model is trained using an objective function stated in terms of margins rather than in terms of $P(y|x)$. However, the model parameters can still be viewed as defining an log-linear or maxent probabilistic representation. CFDs provide a formalism for structured modeling that allows these algorithms and others to be formulated at a level of generality that covers both MRFs of bounded tree width and weighted grammar formalisms like PCFGs.

2 Linear Boolean Models

A *linear Boolean model* (LBM) is a triple $\langle V, F, \Psi \rangle$ where V is a set of Boolean variables with values in $\{0, 1\}$, F is a subset of the set of all assignments to V , and Ψ is an *energy function* $\Psi : V \rightarrow \mathbb{R}$. We call the elements of F the *feasible configurations* of the model. We extend Ψ to configurations $\rho \in F$ with the following “linear” definition:

$$\Psi(\rho) = \sum_{z \in V} \Psi(z)\rho(z) \tag{1}$$

If we view Ψ as a vector in $R^{|V|}$ and ρ as a vector in $\mathbb{B}^{|V|}$ then $\Psi(\rho)$ is simply the inner product of Ψ and ρ . A LBM M defines a probability distribution

$P(\cdot | F, \Psi)$ on feasible configurations $\rho \in F$ as follows.

$$P(\rho | F, \Psi) = \frac{1}{Z(F, \Psi)} e^{-\Psi(\rho)} \quad (2)$$

$$Z(F, \Psi) = \sum_{\rho \in F} e^{-\Psi(\rho)} \quad (3)$$

Given equation (2) we have that an LBM is really just a log-linear or maxent model [11] on a set F under the restrictions that all features are Boolean and that each element of F is uniquely determined by its feature values. A critical issue is how to represent the feasible set F . Before discussing the representation of F , however, we give two examples of representing structured models with LBMs.

3 Markov Random Fields

We first introduce some notation. We consider a set of *variables* V and *domains* $\mathcal{D}(x)$ for each $x \in V$, an *assignment* ρ maps $x \in V$ to $\rho(x) \in \mathcal{D}(x)$; a *partial* assignment σ maps a subset of the variables $\text{dom}(\sigma) \subseteq V$ to appropriate values. We write $\rho' \sqsubseteq \rho$ if $\text{dom}(\rho') \subseteq \text{dom}(\rho)$ and $\rho'(x) = \rho(x) \forall x \in \text{dom}(\rho')$. If ρ is a (possibly partial) assignment on V and $V' \subseteq V$, $\rho|_{V'}$ is the unique assignment such that $\rho|_{V'} \sqsubseteq \rho$ and $\text{dom}(\rho|_{V'}) = \text{dom}(\rho) \cap V'$. If all the variables are *Boolean*, that is $\mathcal{D}(x) = \mathbb{B} = \{0, 1\} \forall x \in V$, the assignment is a *truth* assignment. If ρ is a (possibly partial) assignment, $x \in V$ a variable, and $v \in \mathcal{D}(x)$, $\rho[x := v]$ is the assignment identical to ρ except that $\rho[x := v](x) = v$.

A Markov random field (MRF) consists of variables and energy terms on configurations of those variables. More precisely, we assume a finite set of variables y_1, \dots, y_ℓ with associated domains $\mathcal{D}(y_1), \dots, \mathcal{D}(y_\ell)$. We take the domains $\mathcal{D}(y_j)$ to be finite sets with $|\mathcal{D}(y_i)| \geq 2$. We define a *configuration* to be an assignment ρ of values to the variables. An MRF is a set of such variables plus a set of energy terms Ψ_1, \dots, Ψ_N each of which maps a configuration to a real number. Any such set of energy terms defines a hypergraph on the variables. More specifically, we say that Ψ_k depends on variable y_j if there exists configurations ρ and ρ' which agree on all variables except y_j and such that $\Psi_k(\rho) \neq \Psi_k(\rho')$. Let V_k denote the set of variables on which Ψ_k depends. The sets V_k define a hypergraph on the variables. If $|V_k| = 2$ for all k then these sets define a graph.

An MRF M defines a probability distribution over configurations $P(\rho|M)$

by the following equations:

$$\begin{aligned} P(\rho|M) &= \frac{1}{Z(M)} e^{-\Psi(\rho)} \\ Z(M) &= \sum_{\rho} e^{-\Psi(\rho)} \\ \Psi(\rho) &= \sum_k \Psi_k(\rho) \end{aligned}$$

To represent an MRF as a LBM we must represent a configuration of M as a truth assignment on Boolean variables and represent the energy terms by an energy function on Boolean variables. Given an MRF M we construct Boolean variables of the form " $y_i = v$ " with y_i a variable of M and with $v \in \mathcal{D}(y_i)$. For each energy term Ψ_k with $V_k = \{y_1, \dots, y_m\}$ and each tuple of values v_1, \dots, v_m with $v_i \in \mathcal{D}(y_i)$ we also introduce the Boolean variable " $k, y_1 = v_1 \wedge \dots \wedge y_m = v_m$ ". Of course not all truth assignments to these Boolean variables correspond to configurations of the random field M . In order for a Boolean assignment to be feasible we must have that for each y exactly one of " $y = v_1$ ", \dots , " $y = v_n$ " is true and furthermore " $k, y_1 = v_1 \wedge \dots \wedge y_m = v_m$ " is true if and only if each of " $y_1 = v_1$ ", \dots , " $y_m = v_m$ " is true. Section 5 discusses a method for representing this feasible set of truth assignments. Finally we define the variable energy function as follows.

$$\begin{aligned} \Psi("y = v") &= 0 \\ \Psi("k, y_1 = v_1 \wedge \dots \wedge y_m = v_m") &= \Psi_k(v_1, \dots, v_m) \end{aligned}$$

4 Parse Distributions as LBMs

A CFG in Chomsky normal form is a set of productions of the following form where X , Y and Z are nonterminal symbols and a is a terminal symbol.

$$\begin{aligned} X &\rightarrow YZ \\ X &\rightarrow a \end{aligned}$$

A parse tree is a tree each node of which is labeled by a production of the grammar in the standard way. In a weighted CFG each production $X \rightarrow \gamma$ is assigned an energy (weight) $\Psi(X \rightarrow \gamma)$. For any parse tree y we write $\text{yield}(y)$ for the yield of y , i.e., the sequence of terminal symbols at the leaves of the parse tree. We write $\Psi(y)$ for the total energy of the parse tree y — $\Psi(y)$ is the sum over all nodes of y of the energy of the production used at that node. For a given string x of terminal symbols we have a probability distribution on parse

trees y with $\text{yield}(y) = x$ defined as follows.

$$P(y|x) = \frac{1}{Z(x)} e^{-\Psi(y)} \quad (4)$$

$$Z(x) = \sum_{y: \text{yield}(y)=x} e^{-\Psi(y)} \quad (5)$$

To construct an LBM representation of $P(y|x)$ we first define a set of Boolean variables. Let n be the length of x . First we have a phrase variable " $X_{i,j}$ " for each nonterminal X in the grammar and $1 \leq i < j \leq n+1$. This phrase variable represents the statement that the parse contains a phrase with nonterminal X spanning the string from i to $j-1$ inclusive. Second we have a branch variable " $X_{i,k} \rightarrow Y_{i,j} Z_{j,k}$ " for each production $X \rightarrow YZ$ in the grammar and $1 \leq i < j < k \leq n+1$. A branch variable represents the statement that the parse contains a node labeled with the given production where the left child of the node spans the string from i to $j-1$ and the right child spans j to $k-1$. Finally, we have a terminal variable " $X_{i,i+1} \rightarrow a$ " for each terminal production $X \rightarrow a$ and position i in the input string. A terminal variable represents the statement that the parse tree produces terminal symbol a from nonterminal X at position i . We take V to be the set of all such phrase, branch, and terminal variables. Each parse tree determines a truth assignment to the variables in V and we take F to be the set of assignments corresponding to parse trees. Finally, we must define the energy of each Boolean variable. The variable energy function Ψ is given by the following equations.

$$\begin{aligned} \Psi("X_{i,j}") &= 0 \\ \Psi("X_{i,k} \rightarrow Y_{i,j} Z_{j,k}") &= \Psi(X \rightarrow YZ) \\ \Psi("X_{i,i+1} \rightarrow a") &= \Psi(X \rightarrow a) \end{aligned}$$

5 Case Factor Diagrams (CFDs)

We first introduce some notation. If F is a set of assignments of values to variables then $F[x := v] = \{\rho[x := v] : \rho \in F\}$. If ρ and σ are truth assignments, $\rho \vee \sigma$ is the assignment such that $(\rho \vee \sigma)(x) = 1$ if and only if $\rho(x) = 1$ or $\sigma(x) = 1$. If F_1 and F_2 are sets of truth assignments, $F_1 \vee F_2 = \{\rho \vee \sigma : \rho \in F_1 \text{ and } \sigma \in F_2\}$. The *support* of a truth assignment is the set of variables set to 1 by the assignment.

A *case-factor diagram* represents the feasible set by a search tree over the set of possible truth assignments. The search tree cases on the value of individual variables and factors the feasible set into a product of independent feasible sets when possible. We represent this case-factor search tree by an expression.

Definition 1 A case-factor diagram (CFD) D is an expression generated by the following grammar where x is a Boolean variable; a case expression $\text{case}(x, D_1, D_2)$ must satisfy the constraint that x does not appear in D_1 or D_2 ; and a factor expression $\text{factor}(D_1, D_2)$ must satisfy the constraint that no variable occurs in both D_1 and D_2 .

$$D ::= \text{case}(x, D_1, D_2) \mid \text{factor}(D_1, D_2) \mid \text{unit} \mid \text{empty}$$

We denote by $V(D)$ the set of variables occurring in D .

To define the meaning of CFDs, it is convenient to see all CFD variables as members of a common countably infinite set of variables V . The interpretation $F(D)$ of a CFD D is then a finite set of finite support assignments to V . We use $\bar{0}$ for the totally false assignment (the zero vector). $F(D)$ is defined as follows.

$$\begin{aligned} F(\text{unit}) &= \{\bar{0}\} \\ F(\text{empty}) &= \emptyset \\ F(\text{case}(x, D_1, D_2)) &= F(D_1)[x := 1] \cup F(D_2) \\ F(\text{factor}(D_1, D_2)) &= F(D_1) \vee F(D_2) \end{aligned}$$

Note that, as with ZBDDs, variables that are false in all assignments in $F(D)$ need not be mentioned. In contrast, a BDD must mention any variable required to be false. In contrast to BDDs, ZBDDs give sparse representations of sparse assignments (assignments that are mostly false).

An example consider variables x_1, x_2, \dots and consider the CFD A_i defined as follows.

$$\begin{aligned} A_0 &= \text{unit} \\ A_{i+1} &= \text{case}(x_{i+1}, A_i, A_i) \end{aligned}$$

Under the semantics stated above we have that $F(A_i)$ is the set of all the 2^i truth assignments ρ satisfying the constraint that $\rho(x_j) = 0$ for all $j > i$. As another example, consider B_i defined as follows.

$$\begin{aligned} B_0 &= \text{unit} \\ B_{i+1} &= \text{factor}(\text{case}(x_{i+1}, \text{unit}, \text{unit}), B_i) \end{aligned}$$

We leave it to the reader to verify that $F(B_i) = F(A_i)$. As a third example consider C_i defined as follows.

$$\begin{aligned} C_0 &= \text{unit} \\ C_{i+1} &= \text{case}(x_{i+1}, C_i, \text{empty}) \end{aligned}$$

We have that $F(C_i)$ contains only the single truth assignment ρ such that $\rho(x_j) = 1$ for $j \leq i$ and $\rho(x_j) = 0$ for $j > i$. In general this semantics has

the property that if x does not occur in D then $\rho(x) = 0$ for any assignment $\rho \in F(D)$. Because the two arguments of a **factor** expression cannot share variables, we have that the number of assignments in $F(\mathbf{factor}(D_1, D_2))$ equals the number of assignments in $F(D_1)$ times the number of assignments in $F(D_2)$. We leave it to the reader to verify that any feasible set on any finite set of variables can be represented by a CFD.

The meaning of CFD expressions is independent of their representation as data structures. However, the running time of algorithms depends crucially on that representation. For all the algorithms we discuss, we assume that CFD expressions are represented as *diagrams*, which are DAGs with one node for each distinct subexpression, and edges from the node for an expression to the nodes for its immediate subexpressions. That is, common subexpressions are represented uniquely. For example, the CFD A_i defined above viewed as a tree has 2^i leaves. Viewed as a diagram, however, A_i has only $i + 1$ nodes but 2^i different paths from the root node to the leaf node. The size of a CFD D , denoted $|D|$, is defined to be the number of distinct subexpressions of D (including D itself). In other words, $|D|$ is the number of nodes in the diagram view of D . We will often use the word “node” as a synonym for “expression”. We will also use the standard DAG notions of parent, child, and (directed) path for CFDs.

6 CFDs for MRFs

Here we define a CFD representation of the feasible set for the LBM constructed in section 3. Consider the problem of computing $Z(M)$ for an MRF M . We assume that the variables of M have been given in a fixed order y_1, y_2, \dots, y_n . The assignments to these variables form a tree whose root has a branch for each value of y_1 , the next level branches for each value of y_2 and so on. As variables are assigned, however, the residual hypergraph defined by the energy terms often factors into disjoint sets of terms on disjoint sets of variables. So one can compute $Z(M)$ by factoring the residual problem when possible and, if no factoring is possible, casing out on the value of the next variable (after which more factoring may be possible). This “case-factor process” determines a set of subproblems. The nodes (subexpressions) in the CFD representation of the MRF correspond to the subproblems that arise in this way. Each such subproblem is defined by a subset Σ of the energy terms and a partial assignment ρ to (some of) the variables occurring in Σ .

More formally, consider a subset Σ of the energy terms of M . Let $V(\Sigma)$ be the set of variables on which some energy term in Σ depends, i.e., $V(\Sigma) = \cup_{k \in \Sigma} V_k$. Let ρ be a partial assignment of values to (some of) the variables in $V(\Sigma)$. Note that ρ is defined on the general variables of M rather than the Boolean variables of M' . For each pair of such a subset Σ and partial assignment ρ we now define

a CFD $D(\Sigma, \rho)$. The CFD for the full feasible constraint is $D(\Sigma(M), \emptyset)$ where $\Sigma(M)$ is the set of all energy terms in M and \emptyset is the empty partial assignment. For a given partial assignment ρ we define a graph structure on the energy terms in Σ by saying that there is an edge between two energy terms if there is a variable not assigned a value by ρ on which both terms depend. The key to concise representation is to factor the problem when Σ becomes disconnected. We use the notation $\text{case}(\langle z_1, D_1 \rangle, \langle z_2, D_2 \rangle, \dots, \langle z_m, D_m \rangle)$ as an abbreviation for $\text{case}(z_1, D_1, \text{case}(\langle z_2, D_2 \rangle, \dots, \langle z_n, D_n \rangle))$ where $\text{case}(\langle z, D \rangle)$ is $\text{case}(z, D, \text{empty})$. The CFD $D(\Sigma, \rho)$ is defined as follows.

1. If Σ is disconnected under partial assignment ρ we have

$$D(\Sigma, \rho) = \text{factor}(D(\Sigma_1, \rho|_{V(\Sigma_1)}), \dots, D(\Sigma_n, \rho|_{V(\Sigma_n)}))$$

where $\Sigma_1, \dots, \Sigma_n$ are the connected components of Σ and $\text{factor}(D_1, D_2, \dots, D_n)$ abbreviates $\text{factor}(D_1, \text{factor}(D_2, \dots, D_n))$.

2. Otherwise, if Σ consists of a single constraint Ψ_k and ρ assigns values to all of $V(\Sigma)$, we have the following where $V_k = \{y_1, \dots, y_m\}$ and $v_i = \rho(y_i)$.

$$D(\Sigma, \rho) = \text{case}("k, y_1 = v_1, \dots, y_m = v_m", \text{unit}, \text{empty})$$

3. Otherwise, let y be the earliest variable (under the given variable order) in $V(\Sigma)$ that is not in $\text{dom}(\rho)$. In this case we have the following where $\mathcal{D}(y) = \{v_1, \dots, v_n\}$.

$$D(\Sigma, \rho) = \text{case} \left(\begin{array}{l} \langle "y = v_1", D(\Sigma, \rho[y := v_1]) \rangle, \\ \vdots \\ \langle "y = v_n", D(\Sigma, \rho[y := v_n]) \rangle \end{array} \right)$$

We now show that MRFs of small tree width have concise CFD representations. First we define the notion of tree width.

Definition 2 *We consider a fixed variable order $y_1 \dots y_n$. Consider i with $1 \leq i \leq n + 1$. We say that a variable y_j is past at i if $j < i$ and is future at i if $j \geq i$. We define G_i to be the graph whose nodes are the energy terms of M and where two energy terms are connected by an edge if they both depend on the same future variable (at i). A connected component of G_i (which is a set of energy terms) will be called active if it contains at least one future variable (at i). The tree width of M under the given variable ordering is the maximum over all i of the maximum over all active connected components of G_i of the number of past variables in that component (at i). The tree width of M is the minimum over all orderings of the tree width relative to that ordering.*

Note that an energy term in which all variables are past is not connected to any other energy term. This implies that the inactive components of G_i are singletons in which all variables are past. It also implies that every energy term in an active component must contain at least one future variable.

The above definition of tree width matches standard definitions. Rather than state standard definitions and prove the equivalence we will simply show that under the above definition we have that the tree width of a tree is 1. Suppose that every energy term involves two variables and the graph formed by the energy terms is a tree. Pick a root of the tree and consider a variable ordering that orders parents before children. Now suppose that Σ is an active connected component of G_i . Every energy term in an active Σ must contain a future variable. So in this case the energy terms in an active Σ form a tree every edge of which contains a future variable. In such a tree only the root variable can be past in an ordering that orders parents before children. So the number of past variables in an active component is at most one. We now have the following theorem.

Theorem 1 *Let w be the tree width of M under the given variable ordering. Then $|D(M)|$ is $O(Nd^{w+1})$ where N is the number of energy terms in M and $d = \max_i |\mathcal{D}(y_i)|$.*

Proof: The definition of $D(\Sigma, \rho)$ can be viewed as a set of rules for generating pairs $\langle \Sigma, \rho \rangle$ such that the CFD contains $D(\Sigma, \rho)$. We will call a pair $\langle \Sigma, \rho \rangle$ an anchor pair if there exists an i such that Σ is an active connected component of G_i and ρ assigns values to the past variables of Σ . We assume that $\Sigma(M)$ is a connected component of G_1 (where all variables are future) and that $\Sigma(M)$ contains variables so that $\langle \Sigma(M), \emptyset \rangle$ is an anchor pair. The set of Σ such that there exists an i such that Σ is a connected component of G_i form a tree — as i increases connected components split ultimately terminating in inactive singleton sets. The number of nodes in a tree is not more than twice the number of leaves minus one. Therefore the number of Σ that can appear in anchor pairs is at most $2N - 1$. For a given anchor pair $\langle \Sigma, \rho \rangle$, let $i(\Sigma)$ be the greatest index such that Σ is a connected component of G_i . There are at most w variables in Σ that are past at time $i(\Sigma)$. In any anchor pair $\langle \Sigma, \rho \rangle$ we have that there exists a $j \leq i(\Sigma)$ such that ρ assigns values to the variables in Σ that are before j . The set of ρ satisfying this property forms a tree with at most d^w leaves. Again, since the number of nodes is no larger than twice the number of leaves we have that the number of such assignments ρ is $O(d^w)$. Therefore the total number of anchor pairs in the CFD is $O(Nd^w)$. But each anchor pair generates a certain set of intermediate nodes in the CFD before generating other anchor pairs. The number of intermediate nodes is bounded by the number of triples of the form $\langle \Sigma, \rho[y = v], \Sigma' \rangle$ where $\langle \Sigma, \rho \rangle$ is an anchor pair and Σ' is a component of Σ under $\rho[y = v]$. In each such triple we either have that $\Sigma' = \Sigma$ or we have that Σ and Σ' form an edge in the graph of possible sets Σ . The number of

edges in a tree is no larger than twice the number of leaves. Hence we have that the number of pairs $\langle \Sigma, \Sigma' \rangle$ appearing in these triples is $O(N)$. For a given Σ , the number of assignments of the form $\rho[y = v]$ is $O(d^{w+1})$. So the number of such triples is $O(Nd^{w+1})$. ■

7 CFDs for Parsing

Here we construct a CFD for the feasible set of the LBM defined in Section 4 for a grammar G . We define the CFD $D("X_{i,k}")$ such that the assignments in $F(D("X_{i,k}"))$ are in one-to-one correspondence with the parse trees of the span from i to $k - 1$ with root nonterminal X . The CFD representing the full feasible set of parses is $D("S_{1,n+1}")$. First we define $D("X_{i,k}")$ as follows where $B("X_{i,k}")$ represents the consequences of making " $X_{i,k}$ " true.

$$D("X_{i,k}") = \text{case}("X_{i,k}", B("X_{i,k}"), \text{empty})$$

For $k > i + 1$ we define the consequences $B("X_{i,k}")$ as follows using the multi-branch case notation defined in section 6.

$$B("X_{i,k}") = \text{case}(\langle b_1, B(b_1) \rangle, \dots, \langle b_n, B(b_n) \rangle)$$

where the variables b_p are all possible branch variables of the form " $X_{i,k} \rightarrow Y_{i,j}Z_{j,k}$ ", and $B("X_{i,k} \rightarrow Y_{i,j}Z_{j,k}") = \text{factor}(D("Y_{i,j}"), D("Z_{j,k}"))$.

Finally, if a_i is the i th input symbol, we have

$$B("X_{i,i+1}") = \begin{cases} \text{case}("X_{i,i+1} \rightarrow a_i", \text{unit}, \text{empty}) & \text{if } X \rightarrow a_i \in G \\ \text{empty} & \text{otherwise} \end{cases}$$

This construction has the property that $|D("S_{1,n+1}")|$ is $O(|G|n^3)$ where $|G|$ is the number of productions in the grammar.

8 The Importance of Zero Suppression

We can define the feasible set of parse trees for a given grammar and given input string directly as an MRF on the Boolean variables introduced in section 4. In particular, we can define an MRF with hard constraint energy terms (energy terms that either have infinite energy or zero energy) expressing the following constraints.

1. " $S_{1,n+1}$ " is true.

2. If " $X_{i,k}$ " is true for $k > i + 1$ then there exists " $Y_{i,j}$ " and " $Z_{j,k}$ " where the grammar contains $X \rightarrow YZ$ and " $X_{i,k} \rightarrow Y_{i,j}Z_{j,k}$ " is true.
3. If " $X_{i,k} \rightarrow Y_{i,j}Z_{j,k}$ " is true then " $Y_{i,k}$ " and " $Y_{j,k}$ " are both true.
4. If " $X_{i,i+1}$ " is true then the grammar must contain $X \rightarrow a$ where a is the n th symbol in the input string and " $X_{i,i+1} \rightarrow a$ " is true.
5. If " $X_{i,k} \rightarrow Y_{i,j}Z_{j,k}$ " is true then no other variable of the form " $X_{i,k} \rightarrow W_{i,j'}U_{j',k}$ " is true.
6. If " $Y_{i,j}$ " is true, and is different from " $S_{1,n+1}$ ", then either some variable of the form " $X_{i,k} \rightarrow Y_{i,j}Z_{j,k}$ " is true or some variable of the form " $X_{k,j} \rightarrow Z_{k,i}Y_{i,j}$ " is true.

Constraints 5 and 6 are not implied by 1, 2, 3 and 4. These constraints can be expressed with a SAT problem (a set of disjunctive clauses) where constraint 5 requires $O(G^2n^4)$ clauses. We can then take the resulting MRF and compile it into an **and/or** graph [10] or an algebraic expression [9]. But in this approach both the MRF representation and the compiled form are too large. In the **and/or** graph representation we have **or** nodes representing the choice points corresponding to constraint 2 above. Each branch of an **or** node produces an **and** node. However, without zero suppression (without context-sensitive variable existence), each **and** node must list all the variables that become *false* at that node. For each of the $O(n^3)$ **and** nodes a cubic number of variables become false giving $O(n^6)$ edges in the **and/or** graph.

9 CFDs for Edit Distance

Pair HMMs [12] and weighted finite-state machines [21, 13] have been used to represent trainable weighted edit distance models in text processing and computational biology. As another example of the expressive power of CFDs, we show here how to construct a CFD for the weighted edit distance problem. We will start with a simple context-independent edit cost model, and then indicate how to extend it for context-sensitive edit costs. Consider two strings $a = a_1 \cdots a_m$ and $b = b_1 \cdots b_n$ over a given alphabet V . We view b as being derived from a by insertions, deletions and substitutions. We use the Boolean variable $X_{i,j}$, with $0 \leq i \leq m$ and $0 \leq j \leq n$, to represent the statement that the j -long prefix of ${}_0b_j$ was derived by editing the i -long prefix of ${}_0a_i$. We require the to level statement $X_{m,n}$ to be true — the string b is derived by editing the string a . In addition, we define the following edit variables:

- $A_{i,j,x}$ states that ${}_0b_j$ is derived from ${}_0a_i$ by first deleting the symbol $x \in V$ at position i in a and then deriving ${}_0b_j$ from ${}_0a_{i-1}$.

- $B_{i,j,x}$ states that ${}_0b_j$ is derived from ${}_0a_i$ by first inserting the symbol $x \in V$ at position j in b and then deriving ${}_0b_{j-1}$ from ${}_0a_i$.
- $S_{i,j,x,y}$ states that ${}_0b_j$ is derived from ${}_0a_i$ by substituting $x \in V$ at position j in b for $y \in V$ at position i in a and then deriving ${}_0b_{j-1}$ from ${}_0a_{i-1}$.

Then we define the CFD $D("X_{i,j}")$ as follows:

$$\begin{aligned}
D("X_{i,j}") &= \begin{cases} \text{case}("X_{i,j}", E("X_{i,j}"), \text{empty}) & \text{if } i + j > 0 \\ \text{unit} & \text{otherwise} \end{cases} \\
E("X_{i,j}") &= \text{case}(\langle e_x, E(e_x) \rangle_{x \in V}, \langle f_y, E(f_y) \rangle_{y \in V}, \langle s_{xy}, E(s_{xy}) \rangle_{x \in V, y \in V})
\end{aligned}$$

where $e_x = "A_{i,j,x}"$, $f_y = "B_{i,j,y}"$, and $s_{xy} = "S_{i,j,x,y}"$.

Finally, we define

$$\begin{aligned}
E("A_{i,j,x}") &= \begin{cases} D("X_{i-1,j}") & \text{if } i > 0, x = a_i \\ \text{empty} & \text{otherwise} \end{cases} \\
E("B_{i,j,y}") &= \begin{cases} D("X_{i,j-1}") & \text{if } j > 0, y = b_j \\ \text{empty} & \text{otherwise} \end{cases} \\
E("S_{i,j,x,y}") &= \begin{cases} D("X_{i-1,j-1}") & \text{if } i > 0, j > 0, x = a_i, y = b_j \\ \text{empty} & \text{otherwise} \end{cases}
\end{aligned}$$

It is easy to see that $|D("X_{m,n}")| = O(|V|^2 mn)$, in agreement with the standard dynamic program for computing the best alignment. From this CFD, we can immediately build an LBM for weighted edit distance by setting

$$\begin{aligned}
\Psi("X_{i,j}") &= 0 \\
\Psi("A_{i,j,x}") &= x \text{ deletion cost} \\
\Psi("B_{i,j,y}") &= y \text{ insertion cost} \\
\Psi("S_{i,j,x,y}") &= \text{cost of substituting } x \text{ for } y
\end{aligned}$$

Context-dependent edits can be implemented by subscripting variables with different classes of edit contexts, and enforcing context class constraints appropriate in the CFD.

10 Inference on CFD Models

A *CFD model* $\langle D, \Psi \rangle$ is an LBM whose feasible set is defined by a CFD D and whose energy function Ψ assigns costs to the variables of D . We will now present the main inference algorithms on CFDs.

The Inside Algorithm. We first consider the problem of computing $Z(F(D), \Psi)$ as defined by equation (3). Here we write $Z(D, \Psi)$ as an abbreviated form of $Z(F(D), \Psi)$. It turns out that $Z(D, \Psi)$ can be computed by recursive descent on subexpressions of D using the following equations.

$$\begin{aligned} Z(\text{case}(x, D_1, D_2), \Psi) &= e^{-\Psi(x)} Z(D_1, \Psi) + Z(D_2, \Psi) \\ Z(\text{factor}(D_1, D_2), \Psi) &= Z(D_1, \Psi) Z(D_2, \Psi) \\ Z(\text{unit}, \Psi) &= 1 \\ Z(\text{empty}, \Psi) &= 0 \end{aligned}$$

The correctness of these equations can be proved by induction on the size of D . By caching these computations for each subexpression of D , these equations give a way of computing $Z(D, \Psi)$ in time proportional to $|D|$. These equations are analogous to the inside algorithm used in statistical parsing.

The Viterbi Algorithm. Next we consider the problem of computing minimum energy over the elements of $F(D)$. In particular we define $\Psi^*(D, \Psi)$ as follows.

$$\Psi^*(D, \Psi) = \min_{\rho \in F(D)} \Psi(\rho)$$

We can compute $\Psi^*(D, \Psi)$ using the following equations.

$$\begin{aligned} \Psi^*(\text{case}(z, D_1, D_2), \Psi) &= \min \left(\begin{array}{l} \Psi(z) + \Psi^*(D_1, \Psi), \\ \Psi^*(D_2, \Psi) \end{array} \right) \\ \Psi^*(\text{factor}(D_1, D_2), \Psi) &= \Psi^*(D_1, \Psi) + \Psi^*(D_2, \Psi) \\ \Psi^*(\text{unit}, \Psi) &= 0 \\ \Psi^*(\text{empty}, \Psi) &= +\infty \end{aligned}$$

Again the correctness of these equations can be proved by a direct induction on the size of D . These equations can easily be modified to also compute a truth assignment that achieves the minimum energy. This is a truth assignment of highest probability.

Marginals. Next we consider the problem of computing marginal probabilities of the form $P(z = 1 \mid D, \Psi, \sigma)$ where σ is a partial truth assignment that fixes the values of some of the CFD model variables. We will show that these marginals can be computed in time proportional to $|D| |\text{dom}(\sigma)|$.

The marginal $P(z = 1 \mid D, \Psi, \sigma)$ can be written as follows:

$$\begin{aligned} P(z \mid D, \Psi, \sigma) &= \frac{Z(D, \Psi, \sigma[z := 1])}{Z(D, \Psi, \sigma)} \\ Z(D, \Psi, \sigma) &= \sum_{\rho \in F(D): \sigma \sqsubseteq \rho} e^{-\Psi(\rho)} \end{aligned}$$

So it suffices to be able to compute $Z(D, \Psi, \sigma)$. We now define the auxiliary quantity $Z'(D, \Psi, \sigma) = Z(D, \Psi, \sigma|_{V(D)})$. Our procedure computes $Z(D, \Psi, \sigma)$ by computing $Z'(D', \Psi, \sigma)$ for all subnodes D' of D . Note that the number of such values is $|D|$. The Z' values satisfy the following equations for **factor**, **unit** and **empty** expressions.

$$\begin{aligned} Z'(\text{factor}(D_1, D_2), \Psi, \sigma) &= Z'(D_1, \Psi, \sigma)Z'(D_2, \Psi, \sigma) \\ Z'(\text{unit}, \Psi, \sigma) &= 1 \\ Z'(\text{empty}, \Psi, \sigma) &= 0 \end{aligned}$$

Computing Z' on **case** expressions is more subtle. We now have the following equation where $Z(z, D, D', \Psi, \sigma)$ is defined below.

$$Z'(\text{case}(z, D_1, D_2), \Psi, \sigma) = \begin{cases} e^{-\Psi(z)} Z(z, D, D_1, \Psi, \sigma) & \text{if } \sigma(z) = 1 \\ Z(z, D, D_2, \Psi, \sigma) & \text{if } \sigma(z) = 0 \\ e^{-\Psi(z)} Z(z, D, D_1, \Psi, \sigma) \\ \quad + Z(z, D, D_2, \Psi, \sigma) & \text{otherwise} \end{cases}$$

$Z(z, D, D', \Psi, \sigma)$ expresses the constraint that omitted variables default to 0 in CFDs. If there exists $z' \neq z$ with $\sigma(z') = 1$ where z' occurs in D but not in D' then $Z(z, D, D', \Psi, \sigma) = 0$, otherwise $Z(z, D, D', \Psi, \sigma) = Z'(D', \Psi, \sigma)$.

To analyze the running time of computing $Z(D, \Psi, \sigma)$ we first note that there are a linear number of values needed of the form $Z(z, D', D'', \Psi, \sigma)$. Assuming unit time hash table operations, it is possible to cache the answer to all queries of the form $z \in D'$, for $z' \in \text{dom}(\sigma)$ and D' a node in D , in $O(|D||\sigma|)$ time. Given this cache, each call to $Z(z, D, D', \Psi, \sigma)$ can be computed in time proportional to $|\sigma|$. So the overall computation takes time proportional to $|D||\sigma|$.

The Inside-Outside Algorithm. Using the above conditional probability algorithm to compute $P(z = 1 \mid D, \Psi)$ for all variables z can take $\Omega(|D|^2)$ time. However, a generalization of the inside-outside algorithm can be used to simultaneously compute $P(z = 1 \mid D, \Psi)$ for all variables z in D in $O(|D|)$ time. The value $Z(D, \Psi)$ is the “inside” value associated with D . Intuitively, the outside value of a node in a CFD is the total weight of the “contexts” in which that node appears. We write $D' \preceq D$ to state that node D' occurs in D where we take a node to occur in itself ($D \preceq D$). For a given top level CFD D_{top} and for $D \preceq D_{\text{top}}$ we define the *outside value* $O(D, D_{\text{top}}, \Psi)$ of D (in D_{top}) as follows. First define $O(D_{\text{top}}, D_{\text{top}}, \Psi) = 1$. For $D \neq D_{\text{top}}$ we define $O(D, D_{\text{top}}, \Psi)$ as follows.

$$\begin{aligned}
O(D, D_{\text{top}}, \Psi) &= \sum_{\text{case}(z, D, D') \preceq D_{\text{top}}} O(\text{case}(z, D, D'), D_{\text{top}}, \Psi) e^{-\Psi(z)} \\
&\quad + \sum_{\text{case}(z, D', D') \preceq D_{\text{top}}} O(\text{case}(z, D', D'), D_{\text{top}}, \Psi) \\
&\quad + \sum_{\text{factor}(D', D') \preceq D_{\text{top}}} O(\text{factor}(D, D'), D_{\text{top}}, \Psi) Z(D', \Psi) \\
&\quad + \sum_{\text{factor}(D', D) \preceq D_{\text{top}}} O(\text{factor}(D', D), D_{\text{top}}, \Psi) Z(D', \Psi)
\end{aligned} \tag{6}$$

Once the inside value of every node has been computed, these equations allows the outside values to be computed from the top down, i.e., starting from $O(D_{\text{top}}, D_{\text{top}}, \Psi) = 1$. Note that in this recursion the top level CFD D_{top} does not change. We will write $O(D, \Psi)$ for $O(D, D_{\text{top}}, \Psi)$ when D_{top} is clear from context. Since D_{top} does not change, this top-down calculation can be done in time proportional to the number of nodes. Finally we can compute $P(z = 1 | D_{\text{top}}, \Psi)$ as follows.

Theorem 2

$$\begin{aligned}
P(z = 1 | D_{\text{top}}, \Psi) &= \frac{Z(D_{\text{top}}, \Psi, \emptyset[z := 1])}{Z(D_{\text{top}}, \Psi)} \\
Z(D_{\text{top}}, \Psi, \emptyset[z := 1]) &= \sum_{\text{case}(z, D, D') \preceq D_{\text{top}}} \begin{pmatrix} O(\text{case}(z, D, D'), \Psi) \\ e^{-\Psi(z)} \\ Z(D, \Psi) \end{pmatrix}
\end{aligned}$$

Proof: First we introduce a slight change of notation so as to put the equations in a more standard form for exponential models. Recall that the energy of an assignment ρ is defined as follows.

$$\Psi(\rho) = \sum_x \Psi(x) \rho(x)$$

We can think of ρ as a vector \vec{x} with components x_1, \dots, x_n and we can rewrite $\Psi(\rho)$ as $\Psi(\vec{x})$ as follows.

$$\Psi(\vec{x}) = \sum_{i=1}^n \Psi_i x_i$$

Here we should think of Ψ as a weight vector with components Ψ_i . We can

write $Z(D_{\text{top}}, \Psi)$ as follows.

$$\begin{aligned} Z(D_{\text{top}}, \Psi) &= \sum_{\vec{x} \in F(D_{\text{top}})} e^{-\sum_i \Psi_i x_i} \\ \frac{\partial Z(D_{\text{top}}, \Psi)}{\partial \Psi_i} &= \sum_{\vec{x} \in F(D_{\text{top}})} -x_i e^{-\sum_i \Psi_i x_i} \\ &= -Z(D_{\text{top}}, \Psi, \emptyset[x_i = 1]) \end{aligned}$$

So to compute $Z(D_{\text{top}}, \Psi, \emptyset[x_i = 1])$ it now suffices to compute $\partial Z / \partial \Psi_i$. We now compute $\partial Z / \partial \Psi_i$ by application of the chain rule to the inside rules for calculating $Z(D_{\text{top}}, \Psi)$. This is analogous to the use of the chain rule in computing partial derivatives in backpropagation for tuning the weights of a neural network. We can think of $Z(D_{\text{top}}, \Psi)$ as a function of $Z(D, \Psi)$ — we expand the inside equation for $Z(D', \Psi)$ at every node $D' \neq D$ with $D \preceq D' \preceq D_{\text{top}}$. This function from $Z(D, \Psi)$ to $Z(D_{\text{top}}, \Psi)$ can be differentiated yielding a well defined value for $\partial Z(D_{\text{top}}, \Psi) / \partial Z(D, \Psi)$. We will now show the following.

$$O(D, D_{\text{top}}, \Psi) = \frac{\partial Z(D_{\text{top}}, \Psi)}{\partial Z(D, \Psi)} \quad (7)$$

The proof is by induction on the maximum depth at which D occurs in D_{top} . For the base case we have $D = D_{\text{top}}$ and (7) follows from $\partial Z(D_{\text{top}}, \Psi) / \partial Z(D_{\text{top}}, \Psi) = 1$. For the induction case we can assume (7) for all shallower nodes. We then have the following where W ranges over all the parents of D

$$\frac{\partial Z(D_{\text{top}}, \Psi)}{\partial Z(D, \Psi)} = \sum_W \left(\frac{\partial Z(D_{\text{top}}, \Psi)}{\partial Z(W, \Psi)} \right) \left(\frac{\partial Z(W, \Psi)}{\partial Z(D, \Psi)} \right)$$

Equation (6) lists the four possible types of parents of D . We consider the case where W is $\text{case}(z, D, D')$. In this case we have the following.

$$\begin{aligned} \frac{\partial Z(D_{\text{top}}, \Psi)}{\partial Z(W, \Psi)} &= O(\text{case}(z, D, D'), \Psi) \\ Z(W, \Psi) &= e^{-\Psi(z)} Z(D, \Psi) + Z(D', \Psi) \\ \frac{\partial Z(W, \Psi)}{\partial Z(D, \Psi)} &= e^{-\Psi(z)} \end{aligned}$$

These equations imply the the first line in (6) covers this form of parent W . The second line (6) similarly covers parents of the form $\text{case}(z, D'', D')$. Now suppose that the parent W has the form $\text{factor}(D, D')$. In this case we have the following.

$$\begin{aligned} \frac{\partial Z(D_{\text{top}}, \Psi)}{\partial Z(W, \Psi)} &= O(\text{factor}(D, D'), D) \\ Z(W, \Psi) &= Z(D, \Psi) Z(D', \Psi) \\ \frac{\partial Z(W, \Psi)}{\partial Z(D, \Psi)} &= Z(D', \Psi) \end{aligned}$$

These equations imply that the third line in (6) properly handles parents of the form $\text{factor}(D, D')$. The analysis of the last line in (6) is similar.

Finally, it suffices to show that the right hand side of the second equation in theorem 2 equals $-\partial Z/\partial\Psi(z)$. We have the following where W now ranges over all parents of the variable z .

$$\frac{\partial Z(D_{\text{top}}, \Psi)}{\partial\Psi(z)} = \sum_W \left(\frac{\partial Z(D_{\text{top}}, \Psi)}{\partial Z(W, \Psi)} \right) \left(\frac{\partial Z(W, \Psi)}{\partial\Psi(z)} \right)$$

The parents of the variable z are exactly the nodes of the form $\text{case}(z, D, D')$ for which we have the following.

$$\begin{aligned} Z(\text{case}(z, D, D'), \Psi) &= e^{-\Psi(z)} Z(D, \Psi) + Z(D', \Psi) \\ \left(\frac{\partial Z(\text{case}(z, D, D'), \Psi)}{\partial\Psi(z)} \right) &= -e^{-\Psi(z)} Z(D, \Psi) \end{aligned}$$

These equations imply that the right hand side of the second equation in theorem 2 equals $-\partial Z/\partial\Psi(z)$ as desired. ■

11 Conclusions

We have described a class of structured probabilistic models based on case-factor diagrams. We have also shown that for a given a weighted context free grammar G and input string x the conditional probability $P(y|x)$ can be represented by a CFD model with $O(|G|n^3)$ nodes. We have also shown that any MRF with tree width w in which variables have V possible values and with N energy terms can be represented by a CFD model with $O(NV^w)$ nodes. We have shown that for an arbitrary CFD model, computing the partition function, most likely variable assignment, and the probability of each Boolean variable, can all be done in time linear in the number of nodes. We believe that CFD models will provide a common language for specifying algorithms and stating theorems that can play for structured probabilistic models a similar role to that of BDDs in Boolean inference problems.

Acknowledgments Michael Collins and Fernando Pereira were supported in this work by the National Science Foundation under grants 0347631 and EIA-0205456, respectively.

References

- [1] David Allen and Adnan Darwiche. New advances in inference by recursive conditioning. In *UAI03*, 2003.
- [2] Yasemin Altun and Thomas Hofmann. Large margin methods for label sequence learning. In *8th European Conference on Speech Communication and Technology (EuroSpeech)*, 2003.
- [3] F. Bacchus, S. Dalmao, and T. Pitassi. Algorithms and complexity results for #sat and bayesian inference. In *FOCS*, pages 340–351, 2003.
- [4] F. Bacchus, S. Dalmao, and T. Pitassi. Value elimination: Bayesian inference via backtracking search. In *UAI*, pages 20–28, 2003.
- [5] Craig Boutilier, Nir Friedman, Moises Goldszmidt, and Daphne Koller. Context-specific independence in bayesian networks. In *UAI96*, 1996.
- [6] Randal E. Bryant. Graph-based algorithms for boolean function manipulation. *IEEE Transactions on Computers*, C-35(8):677–691, 1986.
- [7] M. Collins. Parameter estimation for statistical parsing models: Theory and practice of distribution-free methods. In Harry Bunt, John Carroll, and Giorgio Satta, editors, *New Developments in Parsing Technology*. Kluwer, 2004. Revised version of the paper that appeared at IWPT 2001.
- [8] Adnan Darwiche. Recursive conditioning. *Artificial Intelligence*, 125(1-2):5–41, 2001.
- [9] Adnan Darwiche. A differential approach to inference in bayesian networks. *Journal of the ACM*, pages 280–305, May 2003.
- [10] Rina Dechter. And/or search spaces for graphical models. ICS Technical Report, March 2004.
- [11] Stephen Della Pietra, Vincent Della Pietra, and John Lafferty. Inducing features of random fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(4):380–393, 1997.
- [12] R. Durbin, S. Eddy, A. Krogh, and G. Mitchison. *Biological Sequence Analysis: Probabilistic models of proteins and nucleic acids*. Cambridge University Press, 1998.
- [13] Jason Eisner. Parameter estimation for probabilistic finite-state transducers. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, 2002.
- [14] L. Getoor, N. Friedman, D. Koller, and A. Pfeffer. Learning probabilistic relational models. In S. Dzeroski and N. Lavrac, editors, *Relational Data Mining*. Springer-Verlag, 2001.
- [15] Manfred Jaeger. Probabilistic decision graphs: Combining verification and ai techniques for probabilistic inference. *International Journal of Uncertainty Fuzzyness and Knowledge Based Systems*, 12, 2004.
- [16] Keiji Kanazawa, Daphne Koller, and Stuart Russell. Stochastic simulation algorithms for dynamic probabilistic networks. In *UAI95*, 1995.
- [17] John Lafferty, Andrew McCallum, and Fernando Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of ICML-01*, 2001.

- [18] K. L. McMillan. Hierarchical representation of discrete functions, with application to model checking. In *Computer Aided Verification, 6th International Conference*, 1994.
- [19] Shin-ichi Minato. Zero-suppressed BDDs for set manipulation in combinatorial problems. In *Proc. of 30th ACM/IEEE Design Automation Conference (DAC'93)*, pages 272–277. ACM Press, 1993.
- [20] David Poole. Probabilistic Horn abduction and Bayesian networks. *Artificial Intelligence*, 64(1), 1993.
- [21] E. S. Ristad and P. N. Yianilos. Learning string edit distance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(5):522–532, 1998.
- [22] B. Taskar, C. Guestrin, and D. Koller. Max-margin markov networks. In *Neural Information Processing Systems Conference (NIPS03)*, 2003.