# Probabilistic Graphical Models

David Sontag

New York University

Lecture 4, February 21, 2013

# Conditional random fields (CRFs)

- A CRF is a Markov network on variables $\mathbf{X} \cup \mathbf{Y}$, which specifies the conditional distribution

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{c \in C} \phi_c(\mathbf{x}, \mathbf{y}_c)$$

  with partition function

$$Z(\mathbf{x}) = \sum_{\hat{\mathbf{y}}} \prod_{c \in C} \phi_c(\mathbf{x}, \hat{\mathbf{y}}_c).$$

- As before, two variables in the graph are connected with an undirected edge if they appear together in the scope of some factor

- The only difference with a standard Markov network is the normalization term – before marginalized over $\mathbf{X}$ and $\mathbf{Y}$, now only over $\mathbf{Y}$

# Parameterization of CRFs

- We typically parameterize each factor as a log-linear function,

$$\phi_c(\mathbf{x}, \mathbf{y}_c) = \exp\{\mathbf{w}_c \cdot \mathbf{f}_c(\mathbf{x}, \mathbf{y}_c)\}$$

  - $\mathbf{f}_c(\mathbf{x}, \mathbf{y}_c)$ is a feature vector
  - $\mathbf{w}_c$ are weights that are typically learned – we will discuss this extensively in later lectures

- This is without loss of generality: *any* discrete CRF can be parameterized like this (why?)

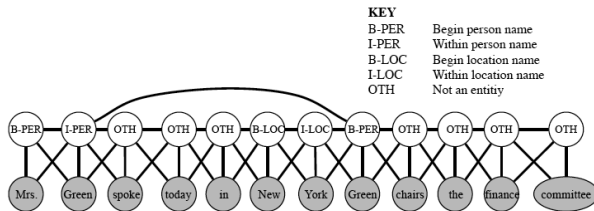- Conditional random fields are in the exponential family:

$$
\begin{aligned}
P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{c \in C} \phi_c(\mathbf{x}, \mathbf{y}_c) &= \exp\left\{ \sum_{c \in C} \mathbf{w}_c \cdot \mathbf{f}_c(\mathbf{x}, \mathbf{y}_c) - \ln Z(\mathbf{w}, \mathbf{x}) \right\} \\
&= \exp\left\{ \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}) - \ln Z(\mathbf{w}, \mathbf{x}) \right\}.
\end{aligned}
$$

# NLP example: named-entity recognition

- Given a sentence, determine the people and organizations involved and the relevant locations:
  "Mrs. Green spoke today in New York. Green chairs the finance committee."
- Entities sometimes span multiple words. Entity of a word not obvious without considering its *context*
- CRF has one variable $X_i$ for each word, which encodes the possible labels of that word
- The labels are, for example, "B-person, I-person, B-location, I-location, B-organization, I-organization"
  - Having beginning (B) and within (I) allows the model to segment adjacent entities

# NLP example: named-entity recognition

The graphical model looks like (called a *skip-chain CRF*):



There are three types of potentials:

- $\phi^1(Y_t, Y_{t+1})$ represents dependencies between neighboring target variables [analogous to transition distribution in a HMM]
- $\phi^2(Y_t, Y_{t'})$ for all pairs $t, t'$ such that $x_t = x_{t'}$, because if a word appears twice, it is likely to be the same entity
- $\phi^3(Y_t, X_1, \cdots, X_T)$ for dependencies between an entity and the word sequence [e.g., may have features taking into consideration capitalization]

**Notice that the graph structure changes depending on the sentence!**

1. Worst-case complexity of probabilistic inference
2. Elimination algorithm
3. Running-time analysis of elimination algorithm (*treewidth*)

# Probabilistic inference

- Today we consider exact inference in graphical models

- In particular, we focus on conditional probability queries,

$$p(\mathbf{Y}|\mathbf{E} = \mathbf{e}) = \frac{p(\mathbf{Y}, \mathbf{e})}{p(\mathbf{e})}$$

(e.g., the probability of a patient having a disease given some observed symptoms)

- Let $\mathbf{W} = \mathcal{X} - \mathbf{Y} - \mathbf{E}$ be the random variables that are neither the query nor the evidence. Each of these joint distributions can be computed by marginalizing over the other variables:

$$p(\mathbf{Y}, \mathbf{e}) = \sum_{\mathbf{w}} p(\mathbf{Y}, \mathbf{e}, \mathbf{w}), \quad p(\mathbf{e}) = \sum_{\mathbf{y}} p(\mathbf{y}, \mathbf{e})$$

- Naively marginalizing over all unobserved variables requires an exponential number of computations

- Does there exist a more efficient algorithm?

# Computational complexity of probabilistic inference

- Here we show that, unless P=NP, there does *not* exist a more efficient algorithm

- We show this by reducing 3-SAT, which is NP-hard, to probabilistic inference in Bayesian networks

- 3-SAT asks about the *satisfiability* of a logical formula defined on $n$ literals $Q_1, \ldots, Q_n$, e.g.

$$(\neg Q_3 \vee \neg Q_2 \vee Q_3) \wedge (Q_2 \vee \neg Q_4 \vee \neg Q_5) \cdots$$

- Each of the disjunction terms is called a *clause*, e.g.

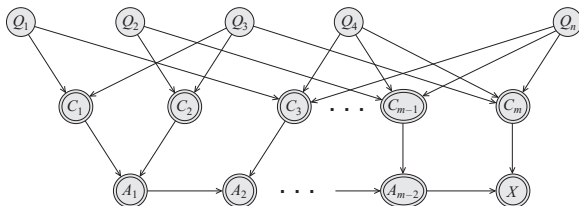$$C_1(q_1, q_2, q_3) = \neg q_3 \vee \neg q_2 \vee q_3$$

In 3-SAT, each clause is defined on at most 3 literals.

- Our reduction also proves that inference in Markov networks is NP-hard (why?)

# Reducing satisfiability to MAP inference

- **Input:** 3-SAT formula with $n$ literals $Q_1, \ldots Q_n$ and $m$ clauses $C_1, \ldots, C_m$



- One variable $Q_i \in \{0, 1\}$ for each literal, $p(Q_i = 1) = 0.5$.
- One variable $C_i \in \{0, 1\}$ for each clause, whose parents are the literals used in the clause. $C_i = 1$ if the clause is satisfied, and 0 otherwise:
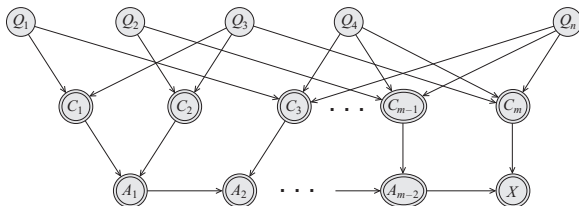
$$p(C_i = 1 \mid \mathbf{q}_{\mathrm{pa}(i)}) = 1[C_i(\mathbf{q}_{\mathrm{pa}(i)})]$$

- Variable $X$ which is 1 if all clauses satisfied, and 0 otherwise:

$$
\begin{aligned}
p(A_i = 1 \mid \mathbf{pa}(A_i)) &= 1[\mathbf{pa}(A_i) = \mathbf{1}], \text{ for } i = 1, \ldots, m - 2 \\
p(X = 1 \mid a_{m-2}, c_m) &= 1[a_{m-2} = 1, c_m = 1]
\end{aligned}
$$

# Reducing satisfiability to MAP inference

- **Input:** 3-SAT formula with $n$ literals $Q_1, \ldots Q_n$ and $m$ clauses $C_1, \ldots, C_m$
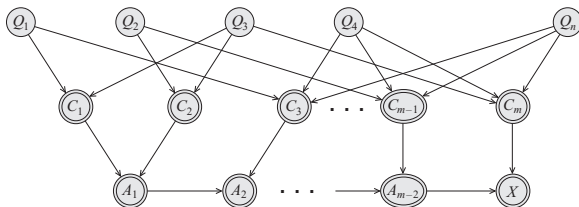


- $p(\mathbf{q}, \mathbf{c}, \mathbf{a}, X = 1) = 0$ for any assignment $\mathbf{q}$ which does not satisfy all clauses
- $p(\mathbf{Q} = \mathbf{q}, \mathbf{C} = \mathbf{1}, \mathbf{A} = \mathbf{1}, X = 1) = \frac{1}{2^n}$ for any satisfying assignment $\mathbf{q}$
- Thus, we can find a satisfying assignment (whenever one exists) by constructing this BN and finding the maximum a posteriori (MAP) assignment:

$$\underset{\mathbf{q},\mathbf{c},\mathbf{a}}{\operatorname{argmax}} \, p(\mathbf{Q} = \mathbf{q}, \mathbf{C} = \mathbf{c}, \mathbf{A} = \mathbf{a} \mid X = 1)$$

- This proves that MAP inference in Bayesian networks and MRFs is NP-hard

# Reducing satisfiability to marginal inference

- **Input:** 3-SAT formula with $n$ literals $Q_1, \ldots Q_n$ and $m$ clauses $C_1, \ldots, C_m$



- $p(X = 1) = \sum_{\mathbf{q,c,a}} p(\mathbf{Q = q, C = c, A = a}, X = 1)$ is equal to the number of satisfying assignments times $\frac{1}{2^n}$

- Thus, $p(X = 1) > 0$ if and only if the formula has a satisfying assignment

- This shows that *marginal inference* is also NP-hard

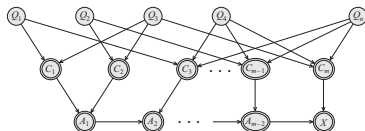# Reducing satisfiability to approximate marginal inference

- Might there exist polynomial-time algorithms that can *approximately* answer marginal queries, i.e. for some $\epsilon$, find $\rho$ such that

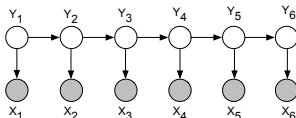$$\rho - \epsilon \leq p(\mathbf{Y} \mid \mathbf{E} = \mathbf{e}) \leq \rho + \epsilon \quad ?$$

- Suppose such an algorithm exists, for any $\epsilon \in (0, \frac{1}{2})$. Consider the following:

  1. Start with $\mathbf{E} = \{ X = 1 \}$
  2. For $i = 1, \ldots, n$:
  3.     Let $q_i = \arg\max_q \; p(Q_i = q \mid \mathbf{E})$
  4.     $\mathbf{E} \leftarrow \mathbf{E} \cup (Q_i = q_i)$

- At termination, $\mathbf{E}$ is a satisfying assignment (if one exists). Pf by induction:
  - In iteration $i$, if $\exists$ satisfying assignment extending $\mathbf{E}$ for **both** $q_i = 0$ *and* $q_i = 1$, then choice in line 3 does not matter
  - Otherwise, suppose $\exists$ satisfying assignment extending $\mathbf{E}$ for $q_i = 1$ but not for $q_i = 0$. Then, $p(Q_i = 1 \mid \mathbf{E}) = 1$ and $p(Q_i = 0 \mid \mathbf{E}) = 0$
  - Even if approximate inference returned $p(Q_i = 1 \mid \mathbf{E}) = 0.501$ and $p(Q_i = 0 \mid \mathbf{E}) = .499$, we would still choose $q_i = 1$

- Thus, it is even NP-hard to *approximately* perform marginal inference!

# Probabilistic inference in practice

- NP-hardness simply says that there **exist** difficult inference problems
- Real-world inference problems are not necessarily as hard as these worst-case instances
- The reduction from SAT created a very complex Bayesian network:



- Some graphs are **easy** to do inference in! For example, inference in hidden Markov models



and other tree-structured graphs can be performed in **linear time**

# Variable elimination (VE)

- Exact algorithm for probabilistic inference in **any** graphical model
- Running time will depend on the *graph structure*
- Uses **dynamic programming** to circumvent enumerating all assignments
- First we introduce the concept for computing marginal probabilities, $p(X_i)$, in Bayesian networks
- After this, we will generalize to MRFs and conditional queries

## Basic idea

- Suppose we have a simple chain, $A \rightarrow B \rightarrow C \rightarrow D$, and we want to compute $p(D)$
- $p(D)$ is a **set** of values, $\{p(D = d), d \in \mathrm{Val}(D)\}$. Algorithm computes sets of values at a time – an entire distribution
- By the chain rule and conditional independence, the joint distribution factors as

$$p(A, B, C, D) = p(A)p(B \mid A)p(C \mid B)p(D \mid C)$$

- In order to compute $p(D)$, we have to marginalize over $A, B, C$:

$$p(D) = \sum_{a,b,c} p(A = a, B = b, C = c, D)$$

# Let's be a bit more explicit…

$$
\begin{array}{llll}
 & P(a^1) & P(b^1 \mid a^1) & P(c^1 \mid b^1) & P(d^1 \mid c^1) \\
+ & P(a^2) & P(b^1 \mid a^2) & P(c^1 \mid b^1) & P(d^1 \mid c^1) \\
+ & P(a^1) & P(b^2 \mid a^1) & P(c^1 \mid b^2) & P(d^1 \mid c^1) \\
+ & P(a^2) & P(b^2 \mid a^2) & P(c^1 \mid b^2) & P(d^1 \mid c^1) \\
+ & P(a^1) & P(b^1 \mid a^1) & P(c^2 \mid b^1) & P(d^1 \mid c^2) \\
+ & P(a^2) & P(b^1 \mid a^2) & P(c^2 \mid b^1) & P(d^1 \mid c^2) \\
+ & P(a^1) & P(b^2 \mid a^1) & P(c^2 \mid b^2) & P(d^1 \mid c^2) \\
+ & P(a^2) & P(b^2 \mid a^2) & P(c^2 \mid b^2) & P(d^1 \mid c^2) \\
\end{array}
$$

$$
\begin{array}{llll}
 & P(a^1) & P(b^1 \mid a^1) & P(c^1 \mid b^1) & P(d^2 \mid c^1) \\
+ & P(a^2) & P(b^1 \mid a^2) & P(c^1 \mid b^1) & P(d^2 \mid c^1) \\
+ & P(a^1) & P(b^2 \mid a^1) & P(c^1 \mid b^2) & P(d^2 \mid c^1) \\
+ & P(a^2) & P(b^2 \mid a^2) & P(c^1 \mid b^2) & P(d^2 \mid c^1) \\
+ & P(a^1) & P(b^1 \mid a^1) & P(c^2 \mid b^1) & P(d^2 \mid c^2) \\
+ & P(a^2) & P(b^1 \mid a^2) & P(c^2 \mid b^1) & P(d^2 \mid c^2) \\
+ & P(a^1) & P(b^2 \mid a^1) & P(c^2 \mid b^2) & P(d^2 \mid c^2) \\
+ & P(a^2) & P(b^2 \mid a^2) & P(c^2 \mid b^2) & P(d^2 \mid c^2) \\
\end{array}
$$

- There is structure to the summation, e.g., repeated $P(c^1|b^1)P(d^1|c^1)$
- Let's modify the computation to first compute

$$P(a^1)P(b^1|a^1) + P(a^2)P(b^1|a^2)$$

## Let's be a bit more explicit...

- Let's modify the computation to first compute

$$P(a^1)P(b^1|a^1) + P(a^2)P(b^1|a^2)$$

and

$$P(a^1)P(b^2|a^1) + P(a^2)P(b^2|a^2)$$

- Then, we get

$$
\begin{array}{llll}
 & (P(a^1)P(b^1 \mid a^1) + P(a^2)P(b^1 \mid a^2)) & P(c^1 \mid b^1) & P(d^1 \mid c^1) \\
+ & (P(a^1)P(b^2 \mid a^1) + P(a^2)P(b^2 \mid a^2)) & P(c^1 \mid b^2) & P(d^1 \mid c^1) \\
+ & (P(a^1)P(b^1 \mid a^1) + P(a^2)P(b^1 \mid a^2)) & P(c^2 \mid b^1) & P(d^1 \mid c^2) \\
+ & (P(a^1)P(b^2 \mid a^1) + P(a^2)P(b^2 \mid a^2)) & P(c^2 \mid b^2) & P(d^1 \mid c^2) \\
\\
 & (P(a^1)P(b^1 \mid a^1) + P(a^2)P(b^1 \mid a^2)) & P(c^1 \mid b^1) & P(d^2 \mid c^1) \\
+ & (P(a^1)P(b^2 \mid a^1) + P(a^2)P(b^2 \mid a^2)) & P(c^1 \mid b^2) & P(d^2 \mid c^1) \\
+ & (P(a^1)P(b^1 \mid a^1) + P(a^2)P(b^1 \mid a^2)) & P(c^2 \mid b^1) & P(d^2 \mid c^2) \\
+ & (P(a^1)P(b^2 \mid a^1) + P(a^2)P(b^2 \mid a^2)) & P(c^2 \mid b^2) & P(d^2 \mid c^2)
\end{array}
$$

- We define $\tau_1 : \mathrm{Val}(B) \to \Re, \quad \tau_1(b^i) = P(a^1)P(b^i|a^1) + P(a^2)P(b^i|a^2)$

# Let's be a bit more explicit...

- We now have

$$
\begin{array}{cccc}
 & \tau_1(b^1) & P(c^1 \mid b^1) & P(d^1 \mid c^1) \\
+ & \tau_1(b^2) & P(c^1 \mid b^2) & P(d^1 \mid c^1) \\
+ & \tau_1(b^1) & P(c^2 \mid b^1) & P(d^1 \mid c^2) \\
+ & \tau_1(b^2) & P(c^2 \mid b^2) & P(d^1 \mid c^2)
\end{array}
$$

$$
\begin{array}{cccc}
 & \tau_1(b^1) & P(c^1 \mid b^1) & P(d^2 \mid c^1) \\
+ & \tau_1(b^2) & P(c^1 \mid b^2) & P(d^2 \mid c^1) \\
+ & \tau_1(b^1) & P(c^2 \mid b^1) & P(d^2 \mid c^2) \\
+ & \tau_1(b^2) & P(c^2 \mid b^2) & P(d^2 \mid c^2)
\end{array}
$$

- We can once more reverse the order of the product and the sum and get

$$
\begin{array}{cc}
 & (\tau_1(b^1)P(c^1 \mid b^1) + \tau_1(b^2)P(c^1 \mid b^2)) \quad P(d^1 \mid c^1) \\
+ & (\tau_1(b^1)P(c^2 \mid b^1) + \tau_1(b^2)P(c^2 \mid b^2)) \quad P(d^1 \mid c^2)
\end{array}
$$

$$
\begin{array}{cc}
 & (\tau_1(b^1)P(c^1 \mid b^1) + \tau_1(b^2)P(c^1 \mid b^2)) \quad P(d^2 \mid c^1) \\
+ & (\tau_1(b^1)P(c^2 \mid b^1) + \tau_1(b^2)P(c^2 \mid b^2)) \quad P(d^2 \mid c^2)
\end{array}
$$

- There are still other repeated computations!

## Let's be a bit more explicit...

- We define $\tau_2 : \mathrm{Val}(C) \to \Re$, with

$$
\begin{array}{rcl}
\tau_2(c^1) & = & \tau_1(b^1)P(c^1|b^1) + \tau_1(b^2)P(c^1|b^2) \\
\tau_2(c^2) & = & \tau_1(b^1)P(c^2|b^1) + \tau_1(b^2)P(c^2|b^2)
\end{array}
$$

- Now we can compute the marginal $p(D)$ as

$$
\begin{array}{ll}
 & \tau_2(c^1) \quad P(d^1 \mid c^1) \\
+ & \tau_2(c^2) \quad P(d^1 \mid c^2)
\end{array}
$$

$$
\begin{array}{ll}
 & \tau_2(c^1) \quad P(d^2 \mid c^1) \\
+ & \tau_2(c^2) \quad P(d^2 \mid c^2)
\end{array}
$$

# What did we just do?

- Our goal was to compute

$$p(D) = \sum_{a,b,c} p(a,b,c,D) = \sum_{a,b,c} p(a)p(b \mid a)p(c \mid b)p(D \mid c)$$

$$= \sum_c \sum_b \sum_a p(D \mid c)p(c \mid b)p(b \mid a)p(a)$$

- We can push the summations inside to obtain:

$$p(D) = \sum_c p(D \mid c) \sum_b p(c \mid b) \underbrace{\sum_a \underbrace{p(b \mid a)p(a)}_{\psi_1(a,b)}}_{\tau_1(b)}$$

- Let's call $\psi_1(A, B) = P(A)P(B|A)$. Then, $\tau_1(B) = \sum_a \psi_1(a, B)$
- Similarly, let $\psi_2(B, C) = \tau_1(B)P(C|B)$. Then, $\tau_2(C) = \sum_b \psi_1(b, C)$
- This procedure is dynamic programming: computation is inside out instead of outside in

## Inference in a chain

- Generalizing the previous example, suppose we have a chain $X_1 \rightarrow X_2 \rightarrow \cdots \rightarrow X_n$ where each variable has $k$ states

- In Problem Set 2, you gave an algorithm to compute $p(X_i)$, for $k = 2$

- For $i = 1$ up to $n - 1$, compute (and cache)

$$p(X_{i+1}) = \sum_{x_i} p(X_{i+1} \mid x_i) p(x_i)$$

- Each update takes $k^2$ time (why?)

- The total running time is $\mathcal{O}(nk^2)$

- In comparison, naively marginalizing over all latent variables has complexity $\mathcal{O}(k^n)$

- We did inference over the joint without ever explicitly constructing it!

# Summary so far

- Worst-case analysis says that marginal inference is NP-hard

- Even approximating it is NP-hard

- In practice, due to the structure of the Bayesian network, we can cache computations that are otherwise computed exponentially many times

- This depends on our having a good **variable elimination ordering**

## Sum-product inference task

- We want to give an algorithm to compute $p(\mathbf{Y})$ for BNs and MRFs

- This can be reduced to the following **sum-product** inference task:

$$\text{Compute} \quad \tau(\mathbf{y}) = \sum_{\mathbf{z}} \prod_{\phi \in \Phi} \phi(\mathbf{z}_{\text{Scope}[\phi] \cap \mathbf{z}}, \mathbf{y}_{\text{Scope}[\phi] \cap \mathbf{Y}}) \quad \forall \mathbf{y},$$

where $\Phi$ is a set of factors or potentials

- For a BN, $\Phi$ is given by the conditional probability distributions for all variables,

$$\Phi = \{\phi_{X_i}\}_{i=1}^n = \{p(X_i \mid \mathbf{X}_{\text{Pa}(X_i)})\}_{i=1}^n,$$

and where we sum over the set $\mathbf{Z} = \mathcal{X} - \mathbf{Y}$

- For Markov networks, the factors $\Phi$ correspond to the set of potentials which we earlier called $C$

  - Sum-product returns an unnormalized distribution, so we divide by $\sum_{\mathbf{y}} \tau(\mathbf{y})$

# Factor marginalization

- Let $\phi(\mathbf{X}, Y)$ be a factor where $\mathbf{X}$ is a set of variables and $Y \notin \mathbf{X}$
- **Factor marginalization** of $\phi$ over $Y$ (also called "summing out $Y$ in $\phi$") gives a new factor:

$$\tau(\mathbf{X}) = \sum_Y \phi(\mathbf{X}, Y)$$

For example,

| | | | |
|---|---|---|---|
| $a^1$ | $b^1$ | $c^1$ | 0.25 |
| $a^1$ | $b^1$ | $c^2$ | 0.35 |
| $a^1$ | $b^2$ | $c^1$ | 0.08 |
| $a^1$ | $b^2$ | $c^2$ | 0.16 |
| $a^2$ | $b^1$ | $c^1$ | 0.05 |
| $a^2$ | $b^1$ | $c^2$ | 0.07 |
| $a^2$ | $b^2$ | $c^1$ | 0 |
| $a^2$ | $b^2$ | $c^2$ | 0 |
| $a^3$ | $b^1$ | $c^1$ | 0.15 |
| $a^3$ | $b^1$ | $c^2$ | 0.21 |
| $a^3$ | $b^2$ | $c^1$ | 0.09 |
| $a^3$ | $b^2$ | $c^2$ | 0.18 |

| | | |
|---|---|---|
| $a^1$ | $c^1$ | 0.33 |
| $a^1$ | $c^2$ | 0.51 |
| $a^2$ | $c^1$ | 0.05 |
| $a^2$ | $c^2$ | 0.07 |
| $a^3$ | $c^1$ | 0.24 |
| $a^3$ | $c^2$ | 0.39 |

# Sum-product variable elimination

- Order the variables **Z** (called the **elimination ordering**)
- Iteratively marginalize out variable $Z_i$, one at a time
- For each $i$,
  1. Multiply all factors that have $Z_i$ in their scope, generating a new product factor
  2. Marginalize this product factor over $Z_i$, generating a smaller factor
  3. Remove the old factors from the set of all factors, and add the new one

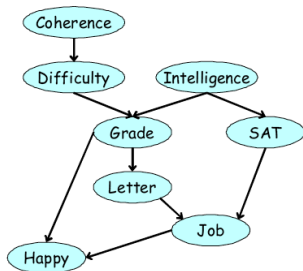**Algorithm 9.1 Sum-Product Variable Elimination algorithm**

    **Procedure** Sum-Product-Variable-Elimination (
      $\Phi$,   // Set of factors
      $\boldsymbol{Z}$,   // Set of variables to be eliminated
      $\prec$   // Ordering on $\boldsymbol{Z}$
    )
1     Let $Z_1, \ldots, Z_k$ be an ordering of $\boldsymbol{Z}$ such that
2       $Z_i \prec Z_j$ iff $i < j$
3     **for** $i = 1, \ldots, k$
4       $\Phi \leftarrow$ Sum-Product-Eliminate-Var$(\Phi, Z_i)$
5     $\phi^* \leftarrow \prod_{\phi \in \Phi} \phi$
6     **return** $\phi^*$

    **Procedure** Sum-Product-Eliminate-Var (
      $\Phi$,   // Set of factors
      $Z$   // Variable to be eliminated
    )
1     $\Phi' \leftarrow \{\phi \in \Phi \; : \; Z \in Scope[\phi]\}$
2     $\Phi'' \leftarrow \Phi - \Phi'$
3     $\psi \leftarrow \prod_{\phi \in \Phi'} \phi$
4     $\tau \leftarrow \sum_Z \psi$
5     **return** $\Phi'' \cup \{\tau\}$

## Example



- What is $p(\mathrm{Job})$? Joint distribution factorizes as:

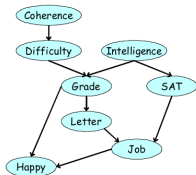$p(C, D, I, G, S, L, H, J) = p(C)p(D|C)p(I)p(G|D, I)p(L|G)P(S|I)P(J|S, L)p(H|J, G)$

with factors

$$\Phi = \{\phi_C(C), \phi_D(C, D), \phi_I(I), \phi_G(G, D, I), \phi_L(L, G),$$
$$\phi_S(S, I), \phi_J(J, S, L), \phi_H(H, J, G)\}$$

- Let's do variable elimination with ordering $\{C, D, I, H, G, S, L\}$ on the board!

# Elimination ordering

- We can pick any order we want, but some orderings introduce factors with much larger scope



| Step | Variable eliminated | Factors used | Variables involved | New factor |
|------|---------------------|--------------|--------------------|------------|
| 1 | $C$ | $\phi_C(C), \phi_D(D,C)$ | $C, D$ | $\tau_1(D)$ |
| 2 | $D$ | $\phi_G(G,I,D), \tau_1(D)$ | $G, I, D$ | $\tau_2(G,I)$ |
| 3 | $I$ | $\phi_I(I), \phi_S(S,I), \tau_2(G,I)$ | $G, S, I$ | $\tau_3(G,S)$ |
| 4 | $H$ | $\phi_H(H,G,J)$ | $H, G, J$ | $\tau_4(G,J)$ |
| 5 | $G$ | $\tau_4(G,J), \tau_3(G,S), \phi_L(L,G)$ | $G, J, L, S$ | $\tau_5(J,L,S)$ |
| 6 | $S$ | $\tau_5(J,L,S), \phi_J(J,L,S)$ | $J, L, S$ | $\tau_6(J,L)$ |
| 7 | $L$ | $\tau_6(J,L)$ | $J, L$ | $\tau_7(J)$ |

- Alternative ordering...

| Step | Variable eliminated | Factors used | Variables involved | New factor |
|------|---------------------|--------------|--------------------|------------|
| 1 | $G$ | $\phi_G(G,I,D), \phi_L(L,G), \phi_H(H,G,J)$ | $G, I, D, L, J, H$ | $\tau_1(I,D,L,J,H)$ |
| 2 | $I$ | $\phi_I(I), \phi_S(S,I), \tau_1(I,D,L,S,J,H)$ | $S, I, D, L, J, H$ | $\tau_2(D,L,S,J,H)$ |
| 3 | $S$ | $\phi_J(J,L,S), \tau_2(D,L,S,J,H)$ | $D, L, S, J, H$ | $\tau_3(D,L,J,H)$ |
| 4 | $L$ | $\tau_3(D,L,J,H)$ | $D, L, J, H$ | $\tau_4(D,J,H)$ |
| 5 | $H$ | $\tau_4(D,J,H)$ | $D, J, H$ | $\tau_5(D,J)$ |
| 6 | $C$ | $\tau_5(D,J), \phi_D(D,C)$ | $D, J, C$ | $\tau_6(D,J)$ |
| 7 | $D$ | $\tau_6(D,J)$ | $D, J$ | $\tau_7(J)$ |

# Choosing an elimination ordering

Set of possible heuristics:

- **Min-neighbors**: The cost of a vertex is the number of neighbors it has in the current graph.

- **Min-weight**: the cost of a vertex is the product of weights (domain cardinality) of its neighbors.

- **Min-fill**: the cost of a vertex is the number of edges that need to be added to the graph due to its elimination.

- **Weighted-Min-Fill**: the cost of a vertex is the sum of weights of the edges that need to be added to the graph due to its elimination. Weight of an edge is the product of weights of its constituent vertices.

Which one better?

- None of these criteria is better than others.

- Often will try several.