

## Lecture 6: Exact inference in Bayes nets. Variable elimination

- What is inference?
- Complexity of exact inference
- Variable elimination

### Examples of MAP queries

- In speech recognition, given a speech signal, one can attempt to reconstruct the most likely sequence of words that could have generated the signal.
- In classification, given the training data and a new example, we want to determine the most probable class label of the new example.

## Queries

Graphical models (directed or undirected) can answer questions about the underlying probability distribution:

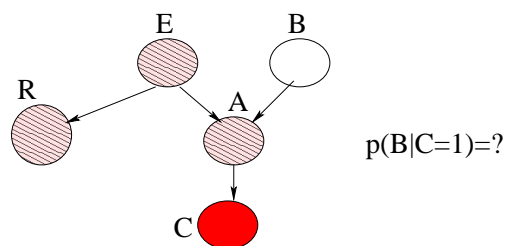
- Conditional or unconditional probability queries:
  - What is the probability of a given value assignment for a subset of variables  $Y$ ?
  - What is the probability of different value assignments for query variables  $Y$  given evidence about variables  $Z$ ? I.e. compute  $p(Y|Z = z)$
- Maximum a posteriori (MAP) queries: given evidence  $Z = z$ , find the most likely assignment of values to the query variables  $Y$ :

$$MAP(Y|Z = z) = \arg \max_y p(Y = y|Z = z)$$

## Complexity of inference

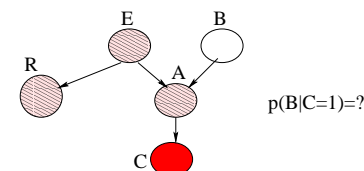
- Given a Bayesian network and a random variable  $X$ , deciding whether  $P(X = x) > 0$  is NP-hard.
- This implies that there is no general inference procedure that will work efficiently for all network configurations
- But for particular families of networks, inference can be done efficiently.
- In other cases, instead of exact inference (computing the probabilities exactly) we will use approximate inference (computing the probabilities with reasonable precision)

### Example of exact inference



$$p(B|C=1) = \frac{p(B, C=1)}{p(C=1)}$$

### Naive solution



$$\begin{aligned} p(B, C=1) &= \sum_{a \in \{0,1\}} \sum_{r \in \{0,1\}} \sum_{e \in \{0,1\}} p(A=a, R=r, E=e, B, C=1) \\ &= \sum_{a,r,e} p(r|e)p(e)p(a|e, B)p(C=1|a) \end{aligned}$$

and same for computing  $p(C=1)$

### A better solution

- Let us re-arrange the sums slightly:

$$\begin{aligned} p(B, C=1) &= \sum_{a,r,e} p(r|e)p(e)p(a|e, B)p(C=1|a) \\ &= \sum_{a,e} p(e)p(a|e, B)p(C=1|a) \sum_r p(r|e) \end{aligned}$$

- Notice that  $\sum_r p(r|e) = 1$ ! But ignore that for the moment. We can call  $\sum_r p(r|e) = m_R(e)$  (because it was obtained by summing out over  $R$  and only depends on  $e$ ).
- Now we have:

$$p(B, C=1) = \sum_a \sum_e p(e)p(a|e, B)p(C=1|a)m_R(e)$$

and we can pick another variable ( $A$  or  $E$ ) to do the same again.

- Instead of  $O(2^n)$  factors, we have to sum over  $O(n \cdot 2^k)$  factors

### Basic idea of variable elimination

- We impose an ordering over the variables, with the query variable coming last
- We maintain a list of “factors”, which depend on given variables
- We sum over the variables in the order in which they appear in the list
- We memorize the result of intermediate computations
- This is a kind of dynamic programming

### A bit of notation

- Let  $X_i$  an evidence variable with observed value  $\hat{x}_i$
- Let the **evidence potential** be an indicator function:

$$\delta(x_i, \hat{x}_i) = 1 \text{ iff } X_i = \hat{x}_i$$

This way, we can turn conditionals into sums as well, e.g.

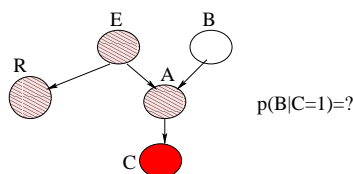
$$p(r|E=1) = \sum_e p(r|e)\delta(e, 1)$$

- This is convenient for notation, but in practice we would take “slices” through the probability tables instead.

### Variable elimination algorithm

- Pick a variable ordering with  $Y$  at the end of the list
- Initialize the active factor list:
  - with the CPDs in a Bayes net
  - with the potentials in a Markov random field
- Introduce the evidence by adding to the active factor list the evidence potentials  $\delta(e, \hat{e})$ , for all the variables in  $E$
- For  $i = 1$  to  $n$ 
  - Take the next variable  $X_i$  from the ordering.
  - Take all the factors that have  $X_i$  as an argument off the active factor list, and multiply them, then sum over all values of  $X_i$ , creating a new factor  $m_{X_i}$
  - Put  $m_{X_i}$  on the active factor list

### Example



- Pick a variable ordering: R, E, C, A, B.
- Initialize the active factor list and introduce the evidence:

$$\text{List: } p(R|E), p(E), p(B), p(A|E, B), p(C|A), \delta(C, 1)$$

- Eliminate R: take  $p(R|E)$  off the list, compute  $m_R(e) = \sum_r p(r|e)$ .

$$\text{List: } p(E), p(B), p(A|E, B), p(C|A), \delta(C, 1), m_R(E)$$

### Example (continued)

- Eliminate E:  $m_E(A, B) = \sum_e p(e)p(a|e, b)m_R(e)$

$$\text{List: } p(B), p(C|A), \delta(C, 1), m_E(A, B)$$

- Eliminate C:  $m_C(a) = \sum_c p(c|a)\delta(C, 1)$

$$\text{List: } p(B), m_E(A, B), m_C(A)$$

- Eliminate A:  $m_A(b) = \sum_a m_E(a, b)m_C(a)$

$$\text{List: } p(B), m_A(B)$$

- The answer we need is a vector with 2 entries:  $p(B=1)m_A(B=1)$  and  $p(B=0)m_A(B=0)$ .

### What about undirected models?

- The algorithm is exactly the same, except that the active factors are initialized with the clique potentials rather than conditional probabilities
- Typically the model has clique potentials associated with nodes,  $\psi(X_i)$ , which makes introduction of evidence very easy:

$$\psi^E(x_i) = \psi(x_i)\delta(x_i, \hat{x}_i)$$

- The normalizing constant almost always cancels out, so the operations are done with unnormalized clique potentials
- The only difference compared to the case of directed models is that usually we do not get factors that are 1 anymore

### Complexity of variable elimination

- We need at most  $O(n)$  multiplications to create one entry in a factor (where  $n$  is the total number of variables)
- If  $m$  is the maximum number of values that a variable can take, a factor depending on  $k$  variables will have  $O(m^k)$  entries
- So it is important to have small factors!
- But the size of the factors depends on the ordering of the variables!
- Choosing an optimal ordering is NP-complete for general networks
- But in special cases a good ordering can be found