Exact Inference on Graphical Model

Reference:

Probabilistic Graphical Model Ch.9, Ch. 10 (Koller & Friedman) CMU, 10-708, Fall 2009 Probabilistic Graphical Models Lectures 8,9,10 (Eric Xing)

Probabilistic Inference

- A Graphical Model specifies a joint distribution P_M(X) over a collection of variables X.
- How can we answer queries/questions about P(X) ? That is, how can we inference using P(X) ?
- Type of queries:
 - 1. Likelihood of evidence/assignments on variables
 - 2. Conditional Probability of some variables (given others).
 - 3. Most Probable Assignment for some variables (given others).

Query 1: Likelihood

Given Evidence E = {X₁=x₁,..., X_D=x_D} specifying some variables' value and let Z={Z₁, ..., Z_k} be variables unspecified , the likelihood of a model M yielding this evidence can be computed by:



Query 1: Likelihood

likelihood of
$$E = \sum_{Z_1} \dots \sum_{Z_K} P_M(Z_1, \dots, Z_K, x_1, \dots, x_D)$$

• This measure is often used as criteria for Model Selection.

Ex. In speech recognition,

Z: words (unspecified) , X: wave sample (specified evidence E)



Query 1: Likelihood

likelihood of E
=
$$P_M(X) = \sum_{Z_1} \dots \sum_{Z_K} P_M(Z_1, \dots, Z_K, x_1, \dots, x_D)$$

Taking special case E = empty , it can also be used to compute Normalizing Const. = Z in MRF as following :

 $(let \widetilde{P}(Z_1...Z_K) = \prod_{cliqueCinM} \phi(C) be unnormalized dist., P(Z_1...Z_K) = \frac{1}{Z} \widetilde{P}(Z_1...Z_K))$

$$\sum_{Z_1} \dots \sum_{Z_K} P(Z_1, \dots, Z_K) = \frac{1}{Z} \sum_{Z_1} \dots \sum_{Z_K} \widetilde{P}(Z_1, \dots, Z_K) = 1$$

==>
$$Z = \sum_{Z_1} ... \sum_{Z_K} \widetilde{P}(Z_1, ..., Z_K)$$

Query 2: Conditional (marginal) Probability

 Given Evidence E = {X₁=x₁,..., X_D=x_D} and some other variables Z={Z₁, ..., Z_k} unspecified , Conditional Probability of Z is given by:

$$P(Z | X) = \frac{P(Z, X)}{P(X)}, \text{ where } P(X) \text{ is given by Query 1}$$

 Sometimes we are interested in only some variables Y in Z, where Z = { Y,W }, then conditional (marginal) prob. of Y is

$$P(Y \mid X) = \sum_{W} P(Z \mid X) = \sum_{W_1} \dots \sum_{W_K} P(Y, W_1 \dots W_K \mid X)$$

Naïve summation over uninterested variables W
yield O(|W|^K) complexity...

Query 2: Conditional (marginal) Probability

Ex. In speech recognition,

Z: words (unspecified), X: wave sample (specified evidence E)

$$P(Z \mid X) = \frac{P(Z, X)}{P(X)}, \text{ where } P(X) \text{ is given by Query 1}$$

A word sequence $Z_1...Z_K$'s prob. given the wave sample $X_1...X_K$

If we only care the 1st word, then:

A person's "Language Model" $P(Z_{t+1}|Z_t)$



Query 3: Most Probable Assignment

 Given Evidence E = {X₁=x₁,..., X_D=x_D} and some other variables Z={Z₁, ..., Z_k} unspecified , Most Probable Assignment of Z is given by:

$$MPA(Z \mid X) = \arg \max_{Z} P(Z \mid X)$$
$$= \arg \max_{Z} \frac{P(X \mid Z)P(Z)}{P(X)} = \arg \max_{Z} P(X \mid Z)P(Z)$$

- MPA is also called "maximum a posteriori configuration" or "MAP inference".
- Note: 1. Even if we have computed Query 2 = P(Z|X), it's intractable to enumerate all possible Z to get argmax_z P(Z|X).



Query 3: Most Probable Assignment

We often just want to "decode words" from the wave sample, That is, we care $Z^* = \operatorname{argmax}_Z P(Z|X)$ but not P(Z|X) itself.

Marginal Maximum:

 $\begin{cases} \arg \max_{Z_1} P(Z_1 \mid X) \\ \arg \max_{Z_2} P(Z_2 \mid X) \\ \arg \max_{Z_3} P(Z_2 \mid X) \end{cases} \implies \text{may give } Z_1 = I', \quad Z_2 = comes', \quad Z_3 = front' \\ \text{(inconsistent decoding)} \end{cases}$

Joint Maximum (MPA) :

 $\underset{Z_{1},Z_{2},Z_{3}}{\operatorname{arg max}} P(Z_{1},Z_{2},Z_{3} \mid X)$ $= > may \ give \ 'I' \ 'come' \ 'from'$ (consistent decoding)



In terms of difficulty, there are 3 types of inference problem.

• Inference which is easily solved with Bayes rule.

Today's focus

• Inference which is tractable using some dynamic programming technique.

(e.g. Variable Elimination or J-tree algorithm)

Inference which is proved intractable
& should be solved using some Approximate Method.
(e.g. Approximation with Optimization or Sampling technique.)



Agenda

- Introduce the concept of "Variable Elimination" in special case of Tree-structured Factor Graph.
- Extend the idea of "VE" to general Factor graph with concept of "Clique Tree".
- See how to extend "VE" to "Most Probable Assignment" (MAP configuration) Problem.



How to get P(E=e) ?

$$P(E=e) = \sum_{A} \sum_{B} \sum_{C} \sum_{D} P(A, B, C, D, E=e)$$

By structure of the BN:

P(A, B, C, D, E) = P(E | D)P(D | C)P(C | B)P(B | A)P(A)

 $P(E) = \sum_{D} \sum_{C} \sum_{B} \sum_{A} P(E \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

We can put summation as right as possible...



How to get P(E=e) ?

$$P(E=e) = \sum_{A} \sum_{B} \sum_{C} \sum_{D} P(A, B, C, D, E=e)$$

By structure of the BN:

P(A, B, C, D, E) = P(E | D)P(D | C)P(C | B)P(B | A)P(A)

 $P(E) = \sum_{D} \sum_{C} \sum_{B} \sum_{A} P(E \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$ $= \sum_{D} P(E \mid D) \sum_{C} P(D \mid C) \sum_{B} P(C \mid B) \sum_{A} P(B \mid A) P(A)$



 $P(E = e) = \sum_{A} \sum_{B} \sum_{C} \sum_{D} P(A)P(B \mid A)P(C \mid B)P(D \mid C)P(E = e \mid D)$ $= \sum_{D} P(E \mid D) \sum_{C} P(D \mid C) \sum_{B} P(C \mid B) \sum_{A} P(B \mid A)P(A)$ F(A,B)A Table size=|A||B|



 $P(E = e) = \sum_{A} \sum_{B} \sum_{C} \sum_{D} P(A)P(B \mid A)P(C \mid B)P(D \mid C)P(E = e \mid D)$ $= \sum_{D} P(E \mid D) \sum_{C} P(D \mid C) \sum_{B} P(C \mid B) \sum_{A} P(B \mid A)P(A)$ $\overline{\Sigma_{A}F(A,B)=M(B)}$ Eliminate "A". A Table size=|B|.



 $P(E = e) = \sum_{A} \sum_{B} \sum_{C} \sum_{D} P(A)P(B \mid A)P(C \mid B)P(D \mid C)P(E = e \mid D)$ = $\sum_{D} P(E \mid D) \sum_{C} P(D \mid C) \sum_{B} P(C \mid B) \sum_{A} P(B \mid A)P(A)$

P(C|B)*M(B) = F(B,C)

A Table size=|B||C|.



 $P(E = e) = \sum_{A} \sum_{B} \sum_{C} \sum_{D} P(A) P(B \mid A) P(C \mid B) P(D \mid C) P(E = e \mid D)$



Eliminate "B". A Table size=|C|.



 $P(E = e) = \sum_{A} \sum_{B} \sum_{C} \sum_{D} P(A)P(B \mid A)P(C \mid B)P(D \mid C)P(E = e \mid D)$ $= \sum_{D} P(E \mid D)\sum_{C} P(D \mid C)\sum_{B} P(C \mid B)\sum_{A} P(B \mid A)P(A)$ $P(D \mid C)M(C) = F(C,D)$ A Table size=|C||D|.







Both Time & Space Complexity are O(|A||B|+|B||C|+|C||D|+|D||E|) → O(|Range|²)

Naïve method complexity = $O(|A||B||C||D||E|) \rightarrow O(|Range|^N)$

$$A - B - C - D - E$$

How about inference on Undirected Model (MRF) ?

$$P(A, B, C, D, E) = \frac{1}{Z}\phi(E, D)\phi(D, C)\phi(C, B)\phi(B, A)$$

$$P(E) = \frac{1}{Z} \sum_{D} \sum_{C} \sum_{B} \sum_{A} \phi(E, D) \phi(D, C) \phi(C, B) \phi(B, A)$$
$$= \frac{1}{Z} \sum_{D} \phi(E, D) \sum_{C} \phi(D, C) \sum_{B} \phi(C, B) \sum_{A} \phi(B, A)$$

The same idea applies !!

$$\begin{array}{c} A \\ f(A,B) \end{array} \\ \begin{array}{c} B \\ f(B,C) \end{array} \\ \begin{array}{c} C \\ f(C,D) \end{array} \\ \begin{array}{c} D \\ f(D,E) \end{array} \\ \begin{array}{c} C \\ f(D,E) \end{array} \\ \begin{array}{c} E \\ f(D,E) \end{array} \\ \begin{array}{c} C \\ f(D,E) \end{array} \\ \end{array} \\ \begin{array}{c} C \\ f(D,E) \end{array} \\ \begin{array}{c} C \\ f(D,E) \end{array} \\ \begin{array}{c} C \\ f(D,E) \end{array} \\ \end{array} \\ \end{array} \\ \begin{array}{c} C \\ f(D,E) \end{array} \\ \end{array} \\ \end{array}$$

From now on, we won't distinguish between BN & MRF. The same algorithm applies to them in a "Factor View".

$$P(A, B, C, D, E) = \frac{1}{Z} \phi(E, D) \phi(D, C) \phi(C, B) \phi(B, A)$$

 $P(A, B, C, D, E) = \mathbf{1}^* P(E \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

All viewed as:
$$\frac{1}{Z}f(E,D)f(D,C)f(C,B)f(B,A)$$



f(A,B)=F(A,B)









M(C)*f(C,D) = F(C,D)

Product !
















If the factor graph is a tree without cycle, then VE can be applied in a similar way.



Eliminate from leaves to root.













Follow the Elimination Process, we can build a "Clique Tree". In which:

- 1. Every node is a F(.) before elimination.
- 2. Every edge is a "message" M(.) passed from F(.) to F(.).





Why is Clique Tree useful ?

If we want to know P(B):





The queried node should get messages from all nodes on the tree to get the marginal distribution.

Why is Clique Tree useful ?

To get marginal distribution of N nodes, we don't need run VE "N times", "2 times" are enough to get all possible messages.



P(B,F) = M(B)f(B,F)M(F)

1st pass: Take a node as root. Run Sum-Product from leaves to root.

2nd pass:

Run Sum-Product from root to leaves.

All marginal dist. can be derived from

- Multiply all M(.) from neighbors by f(.) on this node.
- 2. Eliminate unwanted variables.

Why is Clique Tree useful ?

To get marginal distribution of N nodes, we don't need run VE "N times", "2 times" are enough to get all possible messages.



Complexity :

Elimination on a node F(A,B,C) takes O(|A||B||C|) space & time.

So the algorithm's bottleneck is on elimination for the "Largest Node" on clique tree.

Agenda

- Introduce the concept of "Variable Elimination" in special case of Tree-structured Factor Graph.
- Extend the idea of "VE" to General Factor Graph with concept of "Clique Tree".
- See how to extend "VE" to "Most Probable Assignment" (MAP configuration) Problem.

Some problem ignored earlier: Different "Elimination Orders" have different effect.



Some problem ignored earlier: Different "Elimination Orders" have different effect.

Elimination Order 2: A B C D E



Some problem ignored earlier: Different "Elimination Orders" have different effect.

Elimination Order 2: A B C D E



Some problem ignored earlier: Different "Elimination Orders" have different effect.





Some problem ignored earlier: Different "Elimination Orders" have different effect.

Elimination Order 2: A B C D E



Some problem ignored earlier: Different "Elimination Orders" have different effect.



Some problem ignored earlier: Different "Elimination Orders" have different effect.



Some problem ignored earlier: Different "Elimination Orders" have different effect.

Elimination Order 1: B C D E A



In "Tree" structure factor graph, the optimal "Elimination Order" is just "Elimination from leaves".

If factor graph is not Tree, what's the best elimination order ???

When factor graph is not Tree, we want a Elimination Order "introducing as fewer edges as possible" (then we will have factor size smaller).



Unfortunately, Finding Elimination order with "smallest maximum factor" is NP-hard.

It's fortunate that greedy algorithm works quite well in practical, in which, we just search for the "least-cost" variable to eliminate:

1. If variables have same cardinality

 \rightarrow cost = (# of edges introduced by elimination).

2. If variables have different cardinality

Cost = (# of edges)*(weight by cardinality of node involved)



Decoding 2 person's speech from waves:



Because a factor is a "clique" in undirected representation, we transform Factorial HMM into "undirected" before running VE.



Finding Elimination Order:

Find elimination adding as fewer edges as possible. (greedily)



 $M(Z_1, Y_2, X_1) = \sum_{Y_1} f(Y_1, Z_1) f(Y_1, Y_2) f(Y_1, X_1)$

Finding Elimination Order:

Find elimination adding as fewer edges as possible. (greedily)



 $M(Z_1, Y_2, X_1) = \sum_{Y_1} f(Y_1, Z_1) f(Y_1, Y_2) f(Y_1, X_1)$

Finding Elimination Order:

Find elimination adding as fewer edges as possible. (greedily)



 $M(Y_1, Z_2, X_1) = \sum_{Z_1} f(Y_1, Z_1) f(Z_1, Z_2) f(Z_1, X_1)$

Finding Elimination Order:

Find elimination adding as fewer edges as possible. (greedily)



 $M(Y_1, Z_2, X_1) = \sum_{Z_1} f(Y_1, Z_1) f(Z_1, Z_2) f(Z_1, X_1)$

Finding Elimination Order:



 $M(Y_1, Z_1) = \sum_{X_1} f(Y_1, X_1) f(Z_1, X_1)$

Finding Elimination Order:



 $M(Y_1, Z_1) = \sum_{Y_1} f(Y_1, X_1) f(Z_1, X_1)$ X_1

Finding Elimination Order:



Finding Elimination Order:



Finding Elimination Order:

Find elimination adding as fewer edges as possible. (greedily)



$M(Y_1, Z_2) = \sum_{Z_1} M(Y_1, Z_1) f(Z_1, Z_2)$



Finding Elimination Order:



Finding Elimination Order:


Finding Elimination Order:

Find elimination adding as fewer edges as possible. (greedily)



Finding Elimination Order:

Find elimination adding as fewer edges as possible. (greedily)



Finding Elimination Order:

Find elimination adding as fewer edges as possible. (greedily)



After building a clique tree, we can run "2 passes" on the tree to get all messages M(.) needed for computing marginal.

1st pass:

Take a node as root. Run Sum-Product from leaves to root.

2nd pass:

Run Sum-Product from root to leaves.

All marginal dist. can be derived from

- Multiply all M(.) from neighbors by f(.) on this node.
- 2. Eliminate unwanted variable.



After building a clique tree, we can run "2 passes" on the tree to get all messages M(.) needed for computing marginal.



Example : General Factorial HMM

A clique size=5, intractable most of times.

(No tractable elimination exist...)













Example: A Grid MRF



Generally, we will have clique of "size N" for a N*N grid, which is indeed intractable.

What if some variables $X = \{X_1 ... X_D\}$ are given in Evidence :

Given Evidence { B=b } :



What if some variables $X = \{X_1 ... X_D\}$ are given in Evidence :

Given Evidence { B=b } :



A model with evidence equivalent to another model without evidence.

To infer $P_M(Z|X)$, we transform M to another model M' and infer $P_{M'}(Z)$.

If we can know "which variables will be given", then a intractable model will become a tractable one.

Sometimes we want capture more dependency in a model, which induce intractable inference.

$$P(X,Z) = \frac{1}{Z} f(Z_1, X_1) f(Z_1, X_2) f(Z_1, X_3)$$

$$f(Z_2, X_1) f(Z_2, X_2) f(Z_2, X_3)$$

$$f(Z_3, X_1) f(Z_3, X_2) f(Z_3, X_3)$$

$$f(Z_1, Z_2) f(Z_2, Z_3)$$

$$x_1 \quad x_2 \quad x_3$$

But given X1~X3, we actually run inference on another model M'.

If we can know "which variables will be given", then a intractable model will become a tractable one.

Sometimes we want capture more dependency in a model, which induce intractable inference.

$$P(X = x, Z) = \frac{1}{Z} f_x(Z_1) f_x'(Z_1) f_x''(Z_1)$$

$$f_x(Z_2) f_x'(Z_2) f_x''(Z_2)$$

$$f_x(Z_3) f_x'(Z_3) f_x''(Z_3)$$

$$f(Z_1, Z_2) f(Z_2, Z_3)$$

But given X1~X3, we actually run inference on another model M'.

If we can know "which variables will be given", then a intractable model will become a tractable one.

Sometimes we want capture more dependency in a model, which induce intractable inference.

$$P(X = x, Z) = \frac{1}{Z} f_x'(Z_1, Z_2) f_x'(Z_2, Z_3)$$

But given X1~X3, we actually run inference on another model M'.

Z

 X_2

Z₁

X₁

X₂

If we can know "which variables will be given", then a intractable model will become a tractable one.

Sometimes we want capture more dependency in a model, which induce intractable inference.

$$P(Z | X = x) = \frac{P(X = x, Z)}{P(X = x)} = \frac{\frac{1}{Z} f_x'(Z_1, Z_2) f_x'(Z_2, Z_3)}{\frac{1}{Z} \sum_{Z_1} \sum_{Z_2} \sum_{Z_3} f_x'(Z_1, Z_2) f_x'(Z_2, Z_3)} \begin{bmatrix} z_1 & z_2 & z_3 \\ z_1 & z_2 & z_3 \\ z_1 & z_2 & z_3 \end{bmatrix}$$

But given X1~X3, we actually run inference on another model M'.

If we can know "which variables will be given", then a intractable model will become a tractable one.

Sometimes we want capture more dependency in a model, which induce intractable inference.

$$P(Z \mid X = x) = \frac{P(X = x, Z)}{P(X = x)} = \frac{\frac{1}{Z} f_x'(Z_1, Z_2) f_x'(Z_2, Z_3)}{\frac{1}{Z} \sum_{Z_1} \sum_{Z_2} \sum_{Z_3} f_x'(Z_1, Z_2) f_x'(Z_2, Z_3)} \begin{bmatrix} z_1 & z_2 & z_3 \\ z_1 & z_2 & z_3 \\ z_1 & z_2 & z_3 \end{bmatrix}$$

But given X1~X3, we actually run inference on another model M'.

If we can know "which variables will be given", then a intractable model will become a tractable one.

Sometimes we want capture more dependency in a model, which induce intractable inference.



Even if P(Z|X) can be inferred efficiently, "learning P(X,Z)" is intractable. One solution is model P(Z|X) directly, yielding "CRF" model.

$$P(X,Z) = \frac{1}{Z} f(Z_1, X_1) f(Z_1, X_2) f(Z_1, X_3)$$

$$f(Z_2, X_1) f(Z_2, X_2) f(Z_2, X_3)$$

$$f(Z_3, X_1) f(Z_3, X_2) f(Z_3, X_3)$$

$$f(Z_1, Z_2) f(Z_2, Z_3)$$



Intractable MRF model

$$P(Z \mid X) = \frac{1}{Z'(X)} f_X'(Z_1, Z_2) f_X'(Z_2, Z_3)$$



Tractable CRF Model

Agenda

- Introduce the concept of "Variable Elimination" in special case of Tree-structured Factor Graph.
- Extend the idea of "VE" to General Factor Graph with concept of "Clique Tree".
- See how to extend "VE" to "Most Probable Assignment" (MAP configuration) Problem.

Query 3: Most Probable Assignment

 Given Evidence E = {X₁=x₁,..., X_D=x_D} and some other variables Z={Z₁, ..., Z_k} unspecified , Most Probable Assignment of Z is given by:

$$MPA(Z \mid X) = \arg \max_{Z} P(Z \mid X)$$
$$= \arg \max_{Z} \frac{P(X \mid Z)P(Z)}{P(X)} = \arg \max_{Z} P(X \mid Z)P(Z)$$

$$\underset{Z}{\operatorname{argmax}} P(Z \mid X) \neq \begin{cases} \arg \max_{Z_1} P(Z_1 \mid X) \\ \ldots \\ \arg \max_{Z_K} P(Z_K \mid X) \end{cases}$$

What's the different ?

MPA Goal:

$$\max_{Z} P(Z | X) = \max_{Z_1} \dots \max_{Z_K} P(Z_1 \dots Z_K | X)$$

Likelihood Goal::

(Solved using VE)
$$P(X) = \sum_{Z_1} ... \sum_{Z_K} P(Z_1 ... Z_K, X)$$

Exploring the similarity between "max" & " Σ " is the key to solve MPA using VE.

What's the different ?



 $P(E = e) = \sum_{D} \sum_{C} \sum_{B} \sum_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

$$= \sum_{D} P(E \mid D) \sum_{C} P(D \mid C) \sum_{B} P(C \mid B) \sum_{A} P(B \mid A) P(A)$$

$$M(B) : \text{marginal of B}$$

 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

 $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$ M'(B): ???

$M(B) = max_A F(A,B)$: maxMarginal of B

$$M(B) = \max_{A} F(A, B)$$

F(A,B)	b1	b2	b3
a1	1	3	9
a2	2	5	8
a3	4	7	6

$M(B) = max_A F(A,B)$: maxMarginal of B

 $M(B) = \max_{A} F(A, B)$

F(A,B)	b1	b2	b3	
a1	1	3	9	
a2	2	5	8	
a3	4	7	6	
				I

В	b1	b2	b3
A*(B)	a3		

В	b1	b2	b3
M(B)	4		

$M(B) = \max_{A} F(A,B)$: maxMarginal of B

 $M(B) = \max_{A} F(A, B)$

b1	b2	b3	
1	3	9	
2	5	8	
4	7	6	
	1 2 4	J1 J2 1 3 2 5 4 7	J1 J2 J3 J3 J 1 3 9 3

В	b1	b2	b3
A*(B)	a3	a3	

В	b1	b2	b3
M(B)	4	7	

$M(B) = max_A F(A,B)$: maxMarginal of B

 $M(B) = \max_{A} F(A, B)$

F(A,B)	b1	b2	b3
a1	1	3	9
a2	2	5	8
a3	4	7	6

В	b1	b2	b3
A*(B)	a3	a3	a1

В	b1	b2	b3
M(B)	4	7	9



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$ $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$

F(A,B)

F(A,B)	a1	a2	a3
b1			
b2			
b3			



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

 $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$ $M(B) = \max_{A} F(A,B)$

В	A*(B)	M(B)	
b1	a1		
b2	a3		
b3	a2		

F(A,B)	a1	a2	a3
b1		•••	
b2			
b3	•••		



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$ $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$

F(B,C)=P(C|B)M(B)

P(C B)	b1	b2	b3	В	A*(B)	M(B)
c1				b1	a1	
c2				b2	a3	
c3				b3	a2	
	-					-



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

 $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$

F(B,C)=P(C|B)M(B)

F(B,C)	b1	b2	b3	
c1	•••		•••	
c2	c2			
c3				



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

 $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$ $M(C) = \max_{B} F(B,C)$

С	B*(C)	M(C)	F(B,C)	b1	b2	b3
c1	b3		c1	•••	•••	
c2	b1		c2		•••	
с3	b2		c3	•••		



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

 $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$ F(C,D) = P(D | C)M(C)

P(D C)	c1	c2	c3	_	С	B*(C)
d1					c1	b3
d2					c2	b1
d3				1	c3	b2
	•	•	•			

M(C)

• • •

• • •

...



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

 $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$ F(C,D) = P(D|C)M(C)

F(C,D)	c1	c2	с3	
d1				
d2			•••	
d3				



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

 $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$ $M(D) = \max_{C} F(C,D)$

D	C*(D)	M(D)	F(C,D)	c1	c2	с3
d1	c1		d1			
d2	c2		d2			
d3	c3		d3			
Most Probable Assignment on a Chain



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

 $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$

F(D)=P(E=e|D)M(D)

P(E=e D)	d1	d2	d3	-	D	C*(D)	M(D
е						d1	c1	
					J	d2	c2	
						d3	c3	

Most Probable Assignment on a Chain



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

 $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$

$M = max_{D} F(D)$

D *	Μ	F(D,E=e)	d1	d2	d3
d2		е	•••		•••

What we get ? \rightarrow M = max_{ABCD} P(A,B,C,D,E=e)

What we want ? \rightarrow (A*,B*,C*,D*) = argmax_{ABCD} P(A,B,C,D,E=e)

Most Probable Assignment on a Chain



 $\max_{A,B,C,D} P(A,B,C,D,E=e)$

 $= \max_{D} \max_{C} \max_{B} \max_{A} P(E = e \mid D) P(D \mid C) P(C \mid B) P(B \mid A) P(A)$

 $= \max_{D} P(E = e \mid D) \max_{C} P(D \mid C) \max_{B} P(C \mid B) \max_{A} P(B \mid A) P(A)$ What we want? \rightarrow (A*,B*,C*,D*) = argmax_{ABCD} P(A,B,C,D,E=e) (a1, b1, c2, d2)



Most Probable Assignment on general Graph It's straight forward to generalize algorithm above to case of general graph with similarity of " Σ " and "max".

(The difference is there must be a "traceback" procedure to find the "argmax" after we get "max".)





Summary

- To solve inference problems like "likelihood of X", "P(Z|X)", "Most Probable Assignment", we can use Variable Elimination (e.g. Sum-Product) algorithm
- In case of tree-structured factor graph, we just run "2 passes" of VE from leaves to a root & the reverse.
- In case of general-structured graph, we must find a "good" elimination order inducing smallest "maximum clique", which is often done with greedy method.
- When we know which variables will be given in advance, we can derive much easier model M' from original M with evidence, which is more tractable in Inference & Learning.