

Probabilistic Graph Models:
Factor Graphs, Belief Propagation

ECE/CS 498 DS U/G

Lecture 22: Factor Graphs, Belief Propagation

Ravi K. Iyer

Dept. of Electrical and Computer Engineering

University of Illinois at Urbana Champaign

Announcements

- Course Timeline:
 - Mon 4/13: Factor graphs and belief propagation
 - Wed 4/15: ICA 5; HW4
- MP 3 Checkpoint 1 due **tonight @ 11:59 PM** on Compass 2G
 - .ipynb with Tasks 0-1 completed
 - PDF of slides with answers to Tasks 0-1 (using provided template)
- Final Project
 - Progress report 2 due **Friday April 17 @ 11:59 PM** on Compass2G
 - There should be *substantial* progress with projects by this point (i.e. meaningful results, ML/AI models)
- Discussion section on Friday 4/17
 - Practice with factor graphs and belief propagation

Introduction to Factor Graphs

Hidden Markov Models

Model

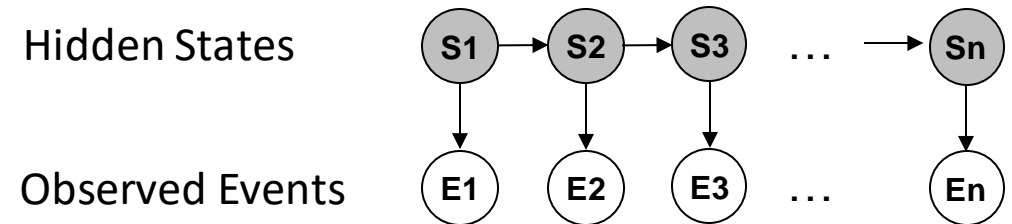
- Set of hidden states $\mathcal{S} = \{\sigma_1, \dots, \sigma_N\}$
- Set of observable events $\mathcal{E} = \{\epsilon_1, \dots, \epsilon_M\}$
- Transition probability matrix A
- Observation matrix B
- Initial distribution of hidden states π

Model assumptions

- An observation depends on its hidden state
- A state variable only depends on the immediate previous state (Markov assumption)
- The future observations and the past observations are **conditionally independent** given the current hidden state

Advantages:

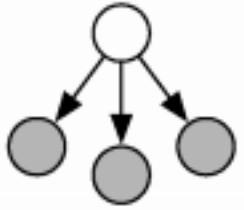
- HMM can model sequential nature of input data (future depends on the past)
- HMM has a linear-chain structure that clearly separates system state and observed events.



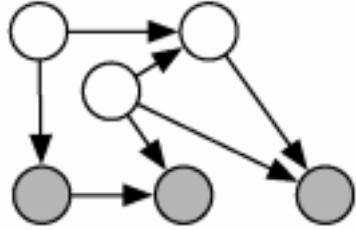
A Hidden Markov model on observed events and system states

$$\begin{aligned} &P(S_1, \dots, S_n, E_1, \dots, E_n) \\ &= P(S_1)P(E_1|S_1) \prod_{i=2}^n P(S_i|S_{i-1})P(E_i|S_i) \end{aligned}$$

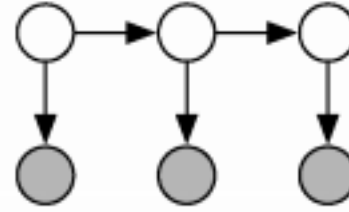
Taxonomy of graphical models



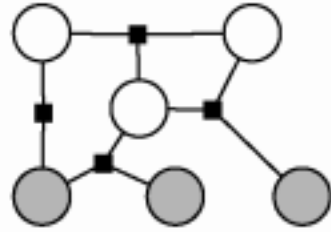
Naïve Bayes



Bayesian Network

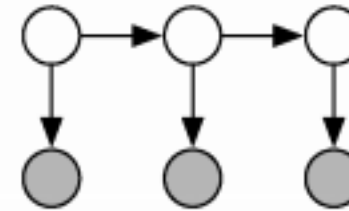
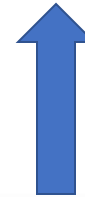


Dynamic Bayesian Network



Factor Graph

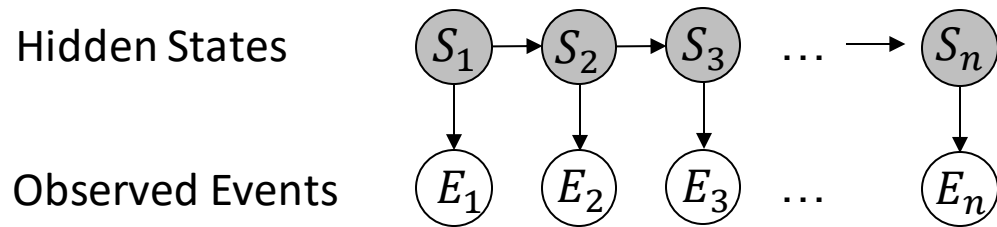
Conditional probabilities and statistical dependencies can be represented by a general type of graph: Factor Graph



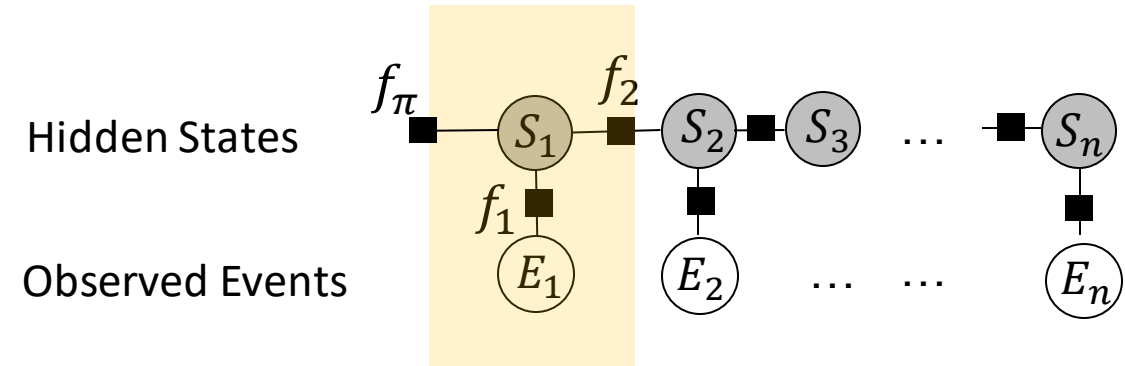
Hidden Markov Model

Conversion of a Hidden Markov Model to a Factor Graph

Hidden Markov Model

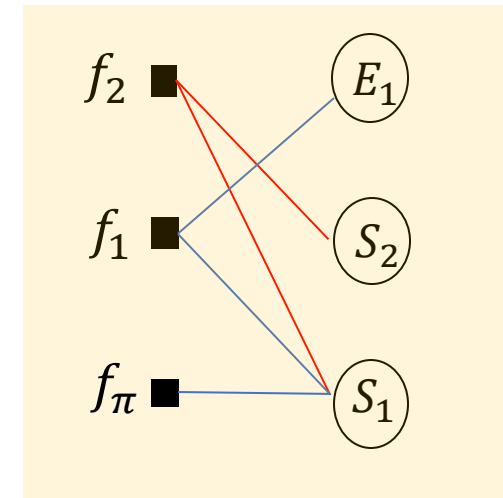


Factor Graph of the HMM



The above **Factor Graph** (FG) is a generalization of the Hidden Markov Model

- Boxes (f_π, f_1, f_2) represents factor function
- In the above case, it maintains the Markov assumption between states



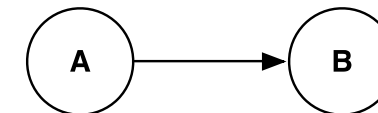
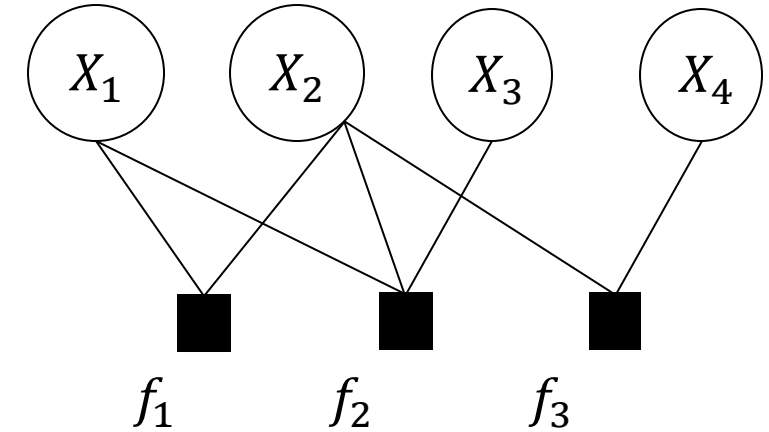
Bipartite graph representation of the FG

Definition of a Factor Graph

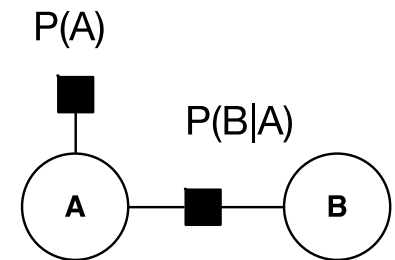
A factor graph is a **bipartite, undirected graph** of **random variables** and **factor functions**.
[Frey et. al. 01].

$G(\text{graph}) = (X, f, E)$; E denotes the edges

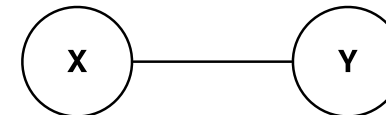
FG can represent both causal and non-causal relations.



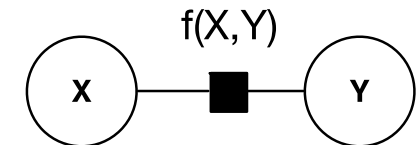
Bayesian Network
(BN)



Factor Graph
equivalent of BN



Undirected Graph



Factor Graph
equivalent of UG

Example Factor function for HMMs

Assume that the state space and observation space are $S = \{\sigma_0, \sigma_1\}$, $E = \{\epsilon_1, \epsilon_2\}$. An example of factor functions is shown.

S	$f_\pi(S)$
σ_0	40
σ_1	25

S_t	E_t	$f_1(S_t, E_t)$
σ_0	ϵ_1	20
σ_0	ϵ_2	15
σ_1	ϵ_1	40
σ_1	ϵ_2	3

S_t	S_{t+1}	$f_2(S_t, S_{t+1})$
σ_0	σ_0	5
σ_0	σ_1	1
σ_1	σ_0	10
σ_1	σ_1	15

- Factor values represents the *affinities* between the related variables
 - E.g., $f_1(\sigma_1, \epsilon_1) > f_1(\sigma_0, \epsilon_1)$ implies that σ_1 and ϵ_1 are more compatible than σ_0 and ϵ_1
- Factor functions don't necessarily represent PDs or joint probability distributions
- How are these values found?
 - Given by expert or from domain knowledge
 - Derived from the data (priors)

Definition of Factor functions

Definition:

- Let \mathbf{D} be a set of random variables. We define a factor f to be a function from $Val(\mathbf{D})$ to \mathbb{R} . A factor is non-negative if all its values are non-negative.
- The set of variables \mathbf{D} is called the scope of the factor f and is denoted as $Scope(f)$.
- $Val(\mathbf{D})$ represents the set of values \mathbf{D} can take.

Example:

A	B	$f(A, B)$
a_0	b_0	30
a_0	b_1	5
a_1	b_0	1
a_1	b_1	10

$$\mathbf{D} = \{A, B\}$$

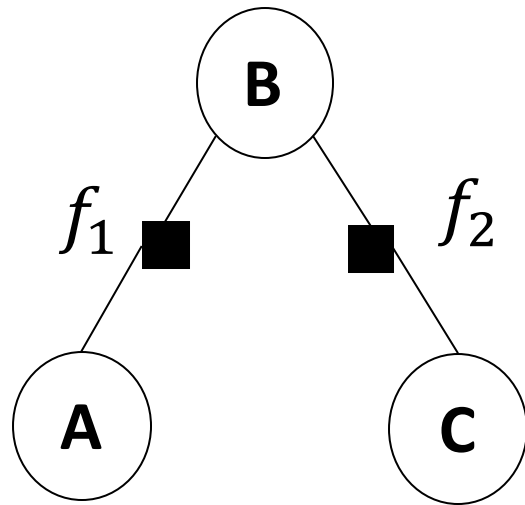
$$Val(\mathbf{D}) = \{(a_0, b_0), (a_0, b_1), (a_1, b_0), (a_1, b_1)\}$$

$$A \in \{a_0, a_1\}$$

$$B \in \{b_0, b_1\}$$

Product of Factor Functions in a Factor Graph

- In HMMs, we derived the joint distribution from the graph representation: $P(S_1, \dots, S_n, E_1, \dots, E_n) = P(S_1)P(E_1|S_1)\prod P(S_i|S_{i-1})P(E_i|S_i)$
- For a Factor Graph, the joint distribution can be derived from the product of factor functions (given that all factor functions are non-negative)



Example Factor Graph
over variables A, B, C .

$$P(A, B, C) = \frac{1}{Z} f_1(A, B) f_2(B, C)$$

where, the normalization Z is given as

$$Z = \sum_{A, B, C} f(A, B, C) = \sum_{A, B, C} f_1(A, B) f_2(B, C)$$

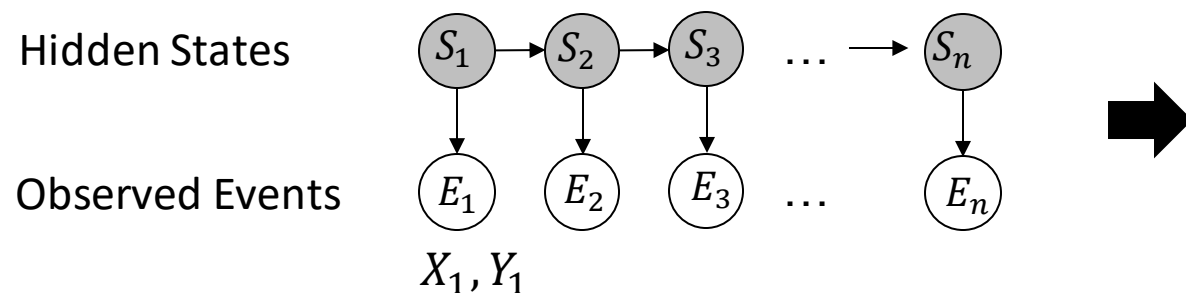
Z is also referred to as the *partition function*.

Conversion of a Hidden Markov Model to a Factor Graph– Two dimension

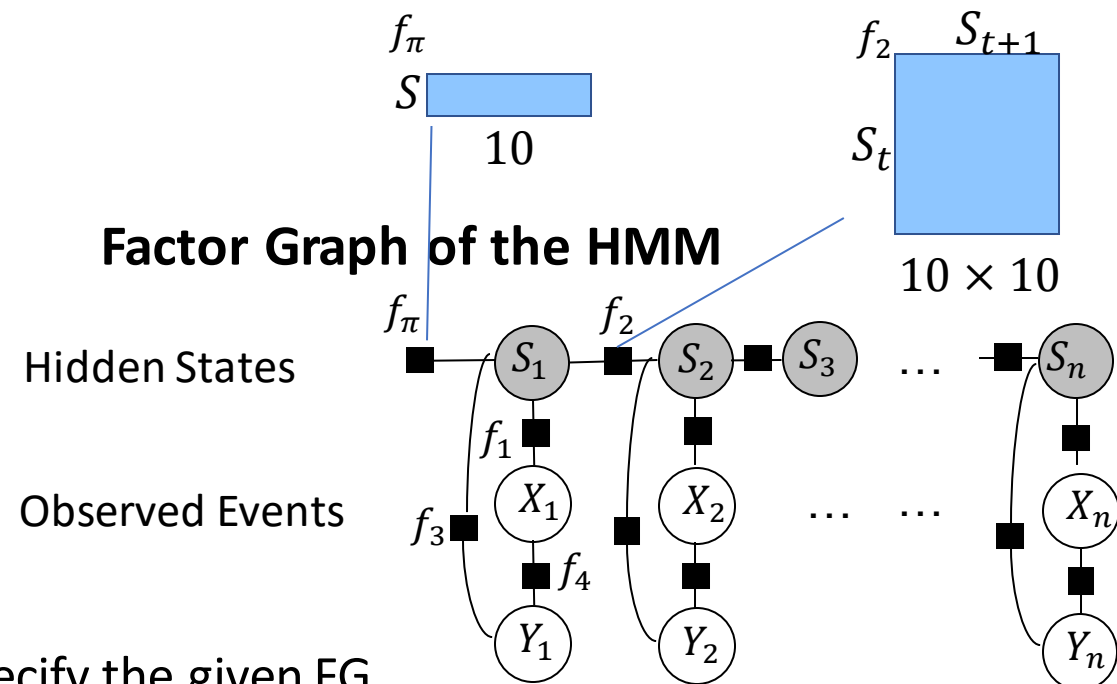
Assume that at each time point, two observations are made corresponding to random variables X and Y .

Example: Let $|S| = 10, |X| = 10, |Y| = 10$

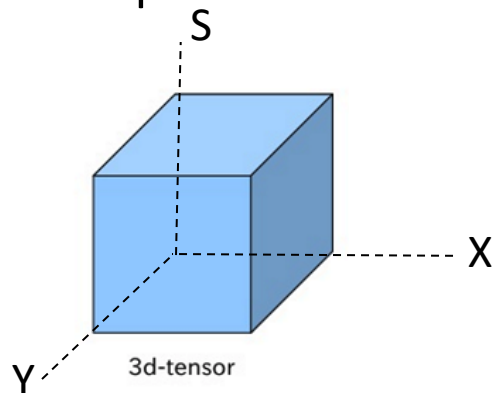
Hidden Markov Model



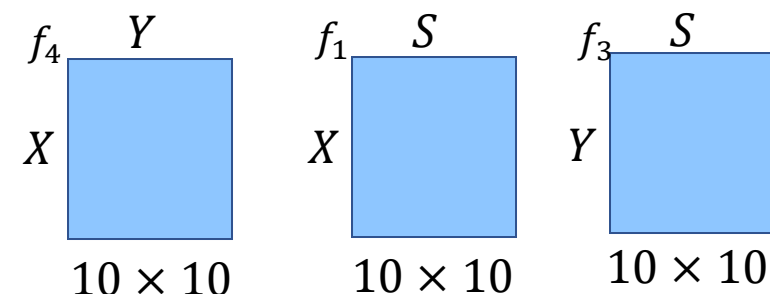
Factor Graph of the HMM



Fewer number of parameters are required are required to specify the given FG.



size of tensor is exponential
 $10 \times 10 \times 10 = 1000$



size of five matrices
 $10 + 10 \times 10 + 10 \times 10 + 10 \times 10 + 10 \times 10 = 410$

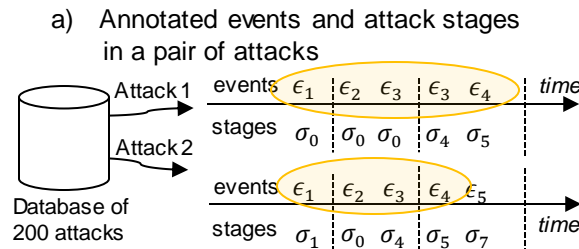
Modeling the credential stealing attack using Factor Graphs - Data

State space of variables

Attack stage: $X = \{\sigma_0, \sigma_1, \dots, \sigma_7\}$

(Observed) Events: $E = \{\epsilon_1, \dots, \epsilon_5\}$

OFFLINE ANNOTATION ON PAST ATTACKS



b) Event-stage annotation table for the attack pair (Attack 1 and Attack 2)

Event	Attack stage
$\{\epsilon_1\}$	$\{\sigma_0, \sigma_1\}$
$\{\epsilon_2\}$	$\{\sigma_0\}$
$\{\epsilon_3\}$	$\{\sigma_4\}$
$\{\epsilon_4\}$	$\{\sigma_5\}$
$\{\epsilon_5\}$	$\{\sigma_7\}$

ϵ_1	vulnerability scan	σ_0	benign
ϵ_2	login	σ_1	discovery
ϵ_3	sensitive_uri	σ_4	privilege escalation
ϵ_4	new_library	σ_5	persistence

- **Attack Information**

- Multi-stage credential stealing attack
- Attack stage $\sigma \in X$ is not observed; however an attack happens in a chain of exploits, thus we have a sequence of events
- Each security event is a known variable ϵ , each takes value from a discrete set of events E

- **Problem statement.** Given a set of security events, infer whether an attack is in progress?

- Goal is to detect and pre-empt the attack

- **Model assumptions**

- There are multivariate relationships among the events
- There is no restriction on order of the relationships (can be non-causal or correlation based)

- Markov Model and Bayesian Networks may not reflect the full range in this scenarios

- Factor graphs can be used for modeling highly complex attacks, where the causal

Modeling the credential stealing attack using Factor Graphs

OFFLINE LEARNING OF FACTOR FUNCTIONS

Example patterns, stages, probabilities, and significance learned from the attack pair

Pattern	Attack stages	Probability in past attacks	Significance (p-value)
$[\epsilon_1, \epsilon_3, \epsilon_4]$	$[\sigma_1, \sigma_4, \sigma_5]$	q_a	p_a
$[\epsilon_1]$	$[\sigma_0 \sigma_1]$	q_b	p_b



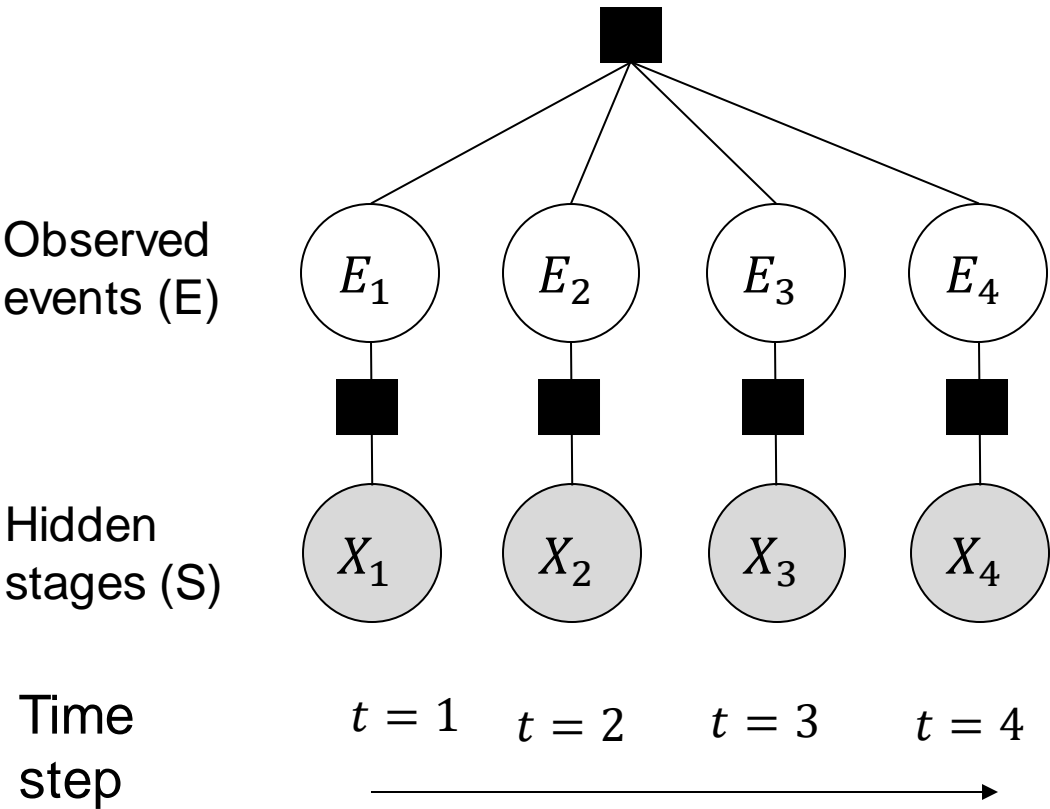
...

■ $f(E) = \exp\{q_E(1 - p_E)\}$

A factor function defined on the learned pattern, stages, and its significance

DETECTION OF UNSEEN ATTACKS

Factor Graph



Advantages and Disadvantages of Factor Graph

Advantage

- Factor graph subsumes HMMs, Markov Random Fields, Bayesian Networks etc.

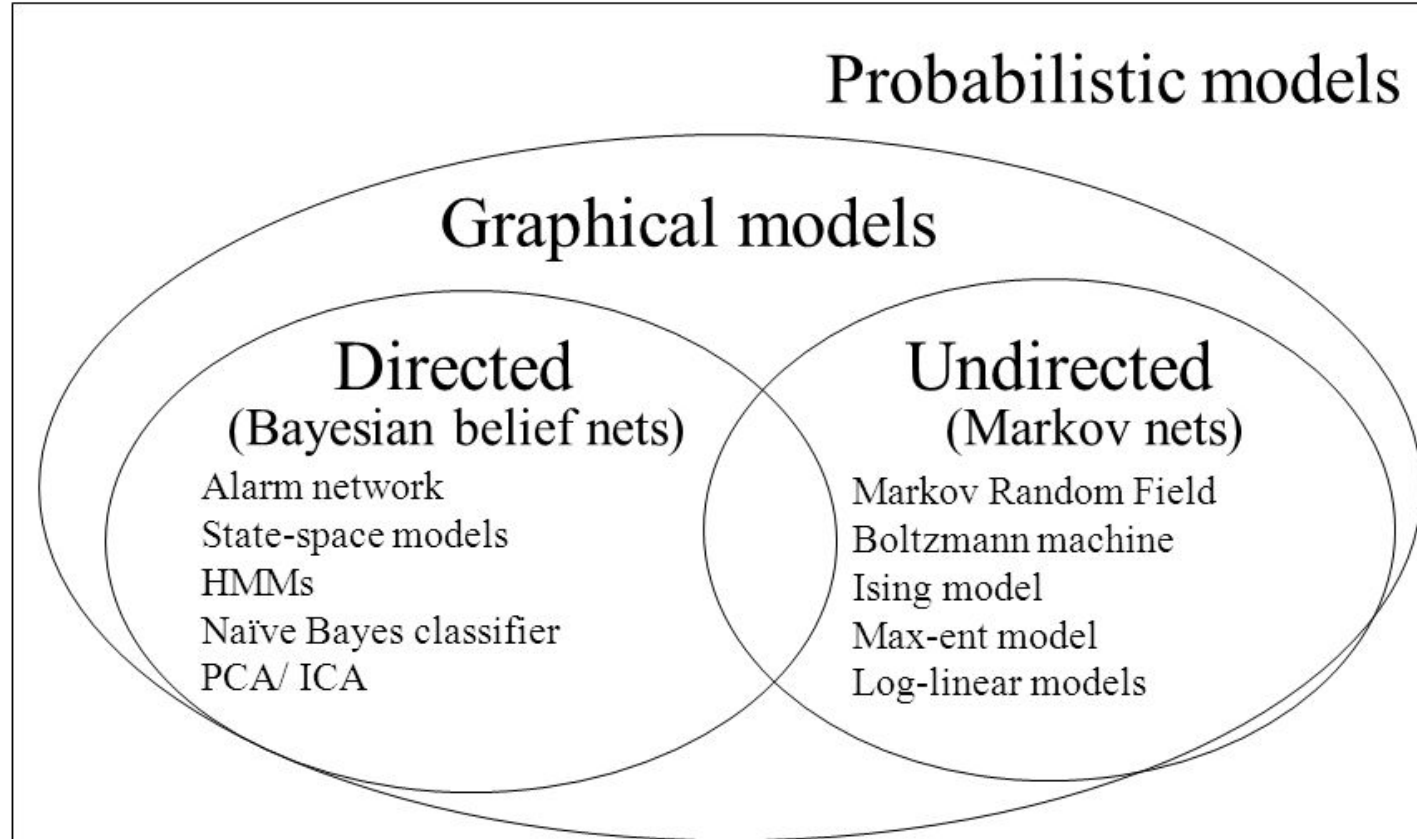
Disadvantage

- Limitations of probabilistic graphs in general

Comment

- If the problem is well represented by specific models such as Bayesian Networks, HMMs, Naïve Bayes or other graphical models then there is no need to generalize your problem as a factor graphs

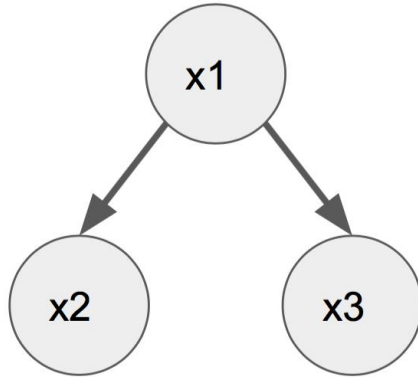
Taxonomy of Graphical Models



Machine Learning, A Probabilistic Perspective, Kevin Murphy, MIT Press

Bayesian Networks vs. Hidden Markov Models vs. Factor Graphs

Bayesian Network

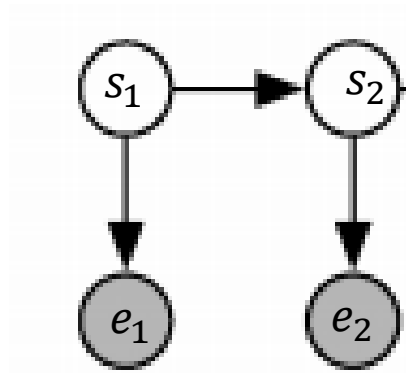


$$p(x_1)p(x_2|x_1)p(x_3|x_1)$$

Product of
conditional
probabilities

Causal relationships

Hidden Markov Model

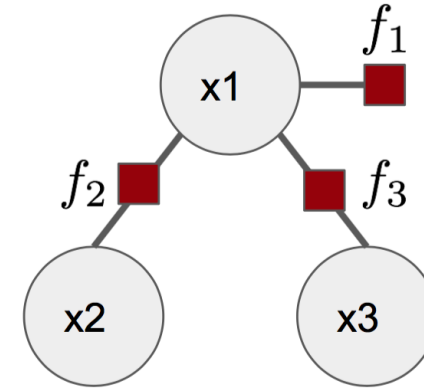


$$p(s_1)p(e_1|s_1)p(s_2|s_1)p(e_2|s_2)$$

Product of
Temporal
dependencies
among variable

Temporal and statistical
dependencies

Factor Graph



$$\frac{1}{Z} f_1(x_1) f_2(x_2, x_1) f_3(x_1, x_3)$$

Product of
dependencies using
univariate, bivariate, or
multivariate functions

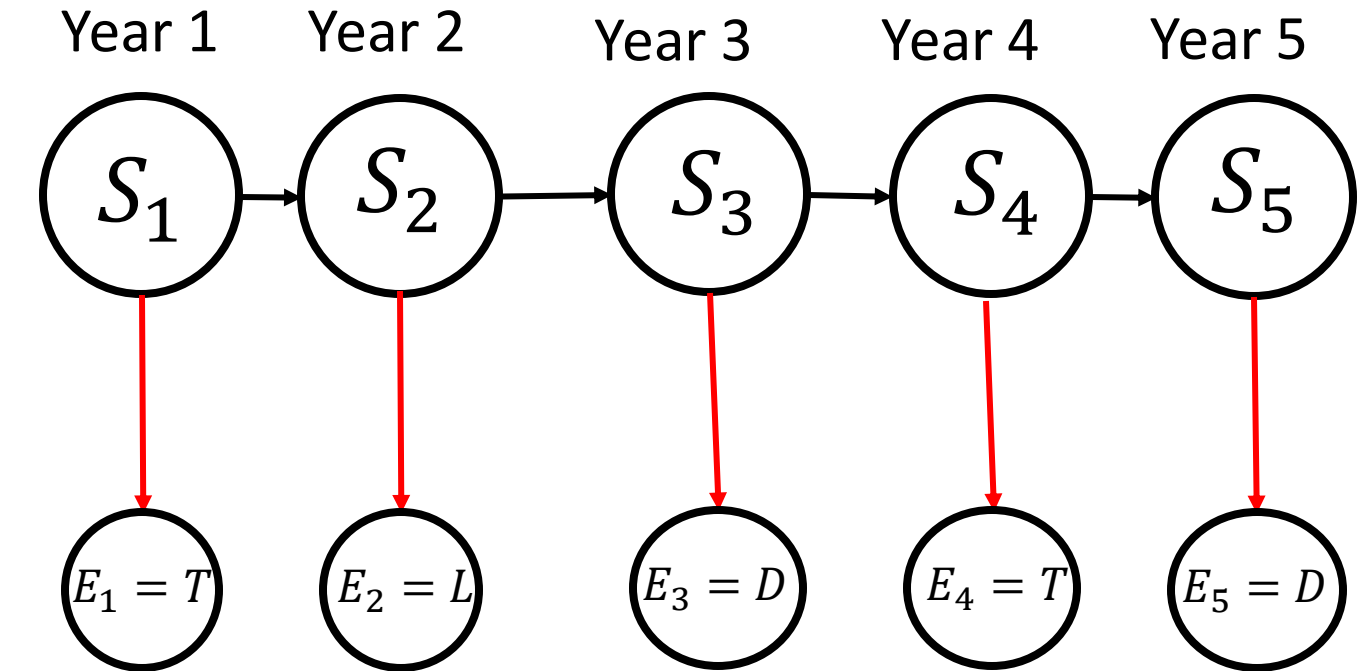
Both types of relations
(including prior on a variable)

Practice with Factor Graph

HMM Example - Paleontological Temperature Model

- State space of hidden states: $S = \{H, C\}$
- State space of observations: $E = \{T, D, L\}$
- Transition probability matrix: A
- Observation Matrix: B
- Initial distribution for the hidden states: π

Given by an oracle



	H	C
H	0.7	0.3
C	0.4	0.6

A

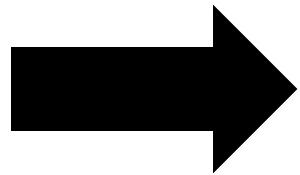
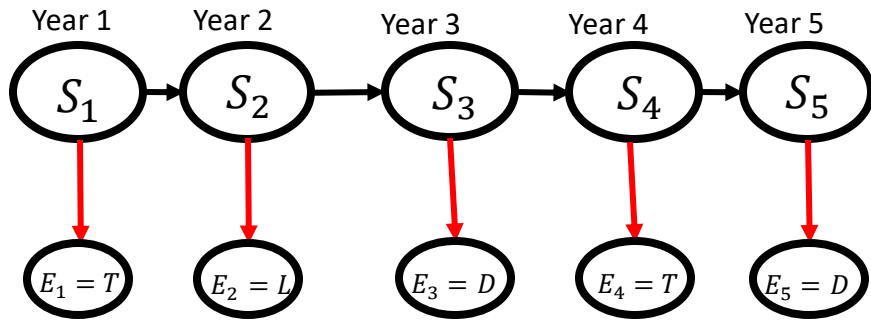
	T	D	L
H	0.1	0.4	0.5
C	0.7	0.2	0.1

B

	H	C
	0.5	0.5

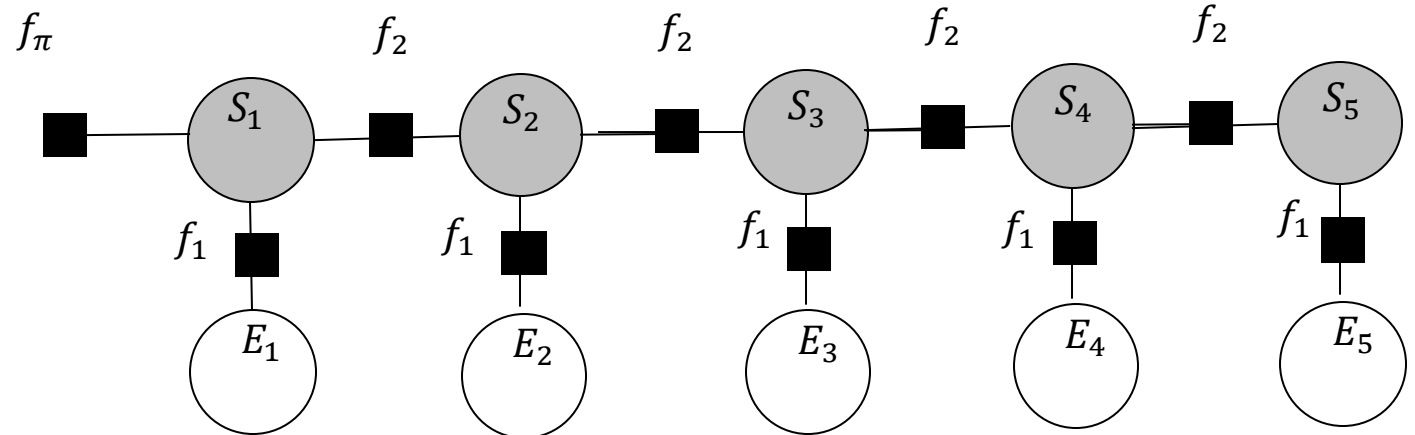
π

First Step – Drawing Factor Graph from HMM



Hidden States

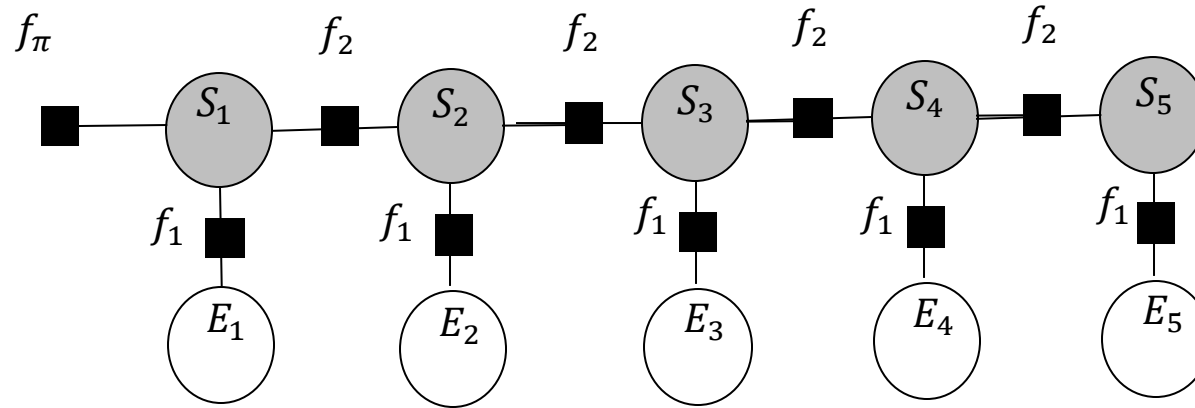
Observed Events



Why are the factor functions...

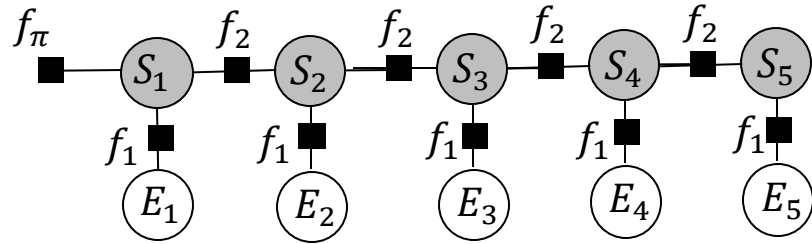
- Between every pair of states the same?
- Between every pair of state and observation the same?

Next – Figuring out the Factor Functions



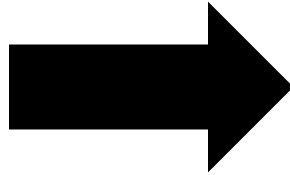
- f_π needs to capture the prior probabilities for the states
- f_1 needs to capture the affinity between observations and states
- f_2 needs to capture the affinity between consecutive states

Next – Figuring out the Factor Functions



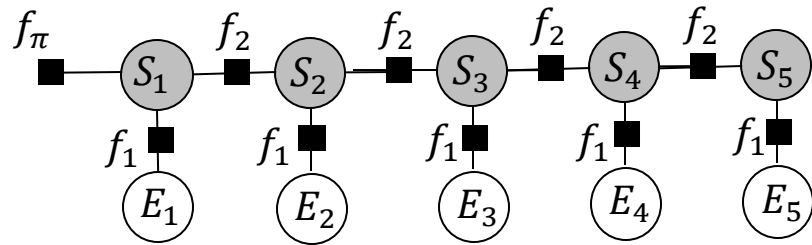
- f_π needs to capture the prior probabilities for the states

$$\begin{array}{cc} H & C \\ [0.5 & 0.5] \\ \pi \end{array}$$



S_1	$f_\pi(S_1)$
H	0.5
C	0.5

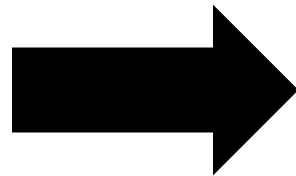
Next – Figuring out the Factor Functions



- f_1 needs to capture the affinity between observations and states. (i.e., $P(E_i|S_i)$)

	T	D	L
H	0.1	0.4	0.5
C	0.7	0.2	0.1

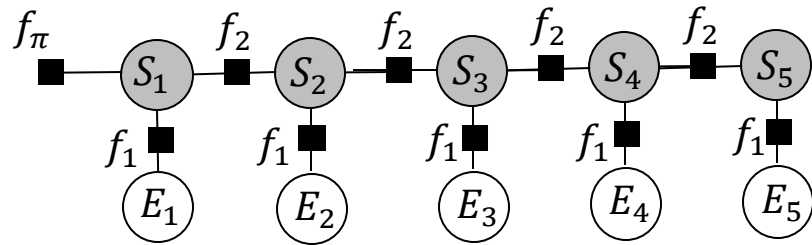
B



S_i	E_i	$f_1(S_i, E_i)$
H	T	0.1
	D	0.4
	L	0.5
C	T	0.7
	D	0.2
	L	0.1

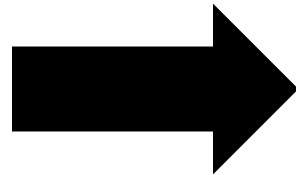
**Note that the values in this table don't sum to 1
 $\Rightarrow f_1$ is not a joint probability but a conditional probability!**

Next – Figuring out the Factor Functions



- f_2 needs to capture the affinity between consecutive states. (i.e., $P(S_{i+1}|S_i)$)

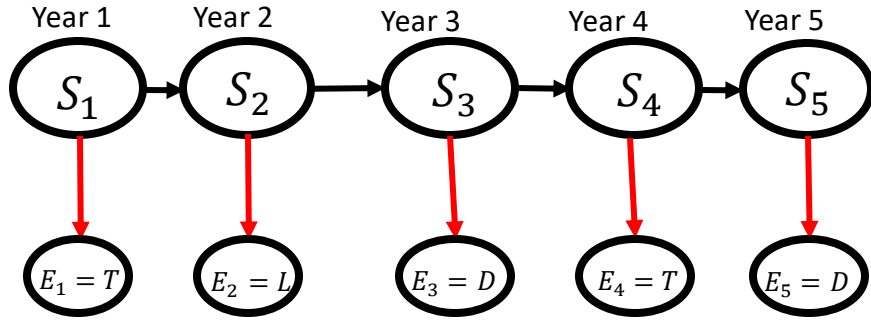
$$\begin{array}{cc} & \begin{matrix} H & C \end{matrix} \\ \begin{matrix} H \\ C \end{matrix} & \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix} \\ & A \end{array}$$



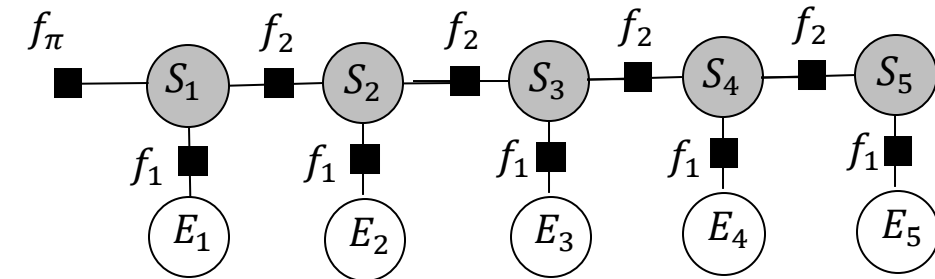
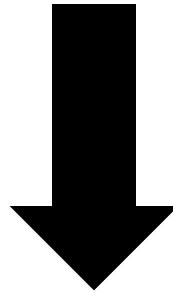
S_i	S_{i+1}	$f_2(S_i, S_{i+1})$
H	H	0.7
	C	0.3
C	H	0.4
	C	0.6

**Note that the values in this table don't sum to 1
 $\Rightarrow f_2$ is not a joint probability but a conditional probability!**

After That – Calculating the Joint



$$P(S_1, \dots, S_5, E_1, \dots, E_5) = P(S_1)P(E_1|S_1) \prod_{i=2}^5 P(S_i|S_{i-1})P(E_i|S_i)$$



$$P(S_1, \dots, S_5, E_1, \dots, E_5) =$$

$$\frac{1}{Z} f_\pi(S_1) f_1(S_1, E_1) \prod_{i=2}^5 f_2(S_{i-1}, S_i) f_1(S_i, E_i)$$

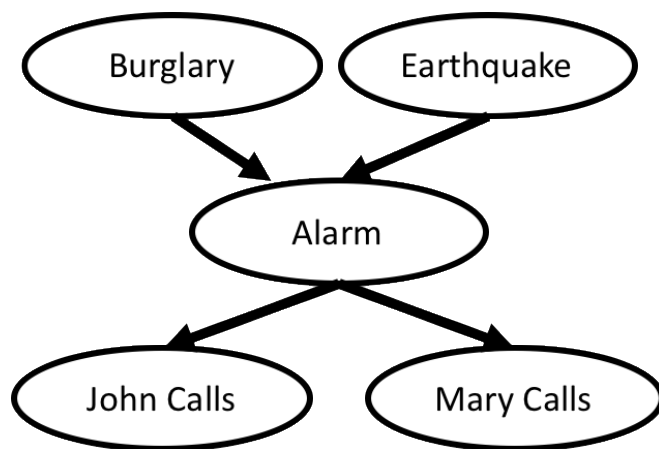
$$Z = \sum_{S_i \in \{H, C\}, E_i \in \{T, D, L\}} f_\pi(S_1) f_1(S_1, E_1) \prod_{i=2}^5 f_2(S_{i-1}, S_i) f_1(S_i, E_i)$$

Belief Propagation

Inference on Graphical Models

Problems we have already looked at:

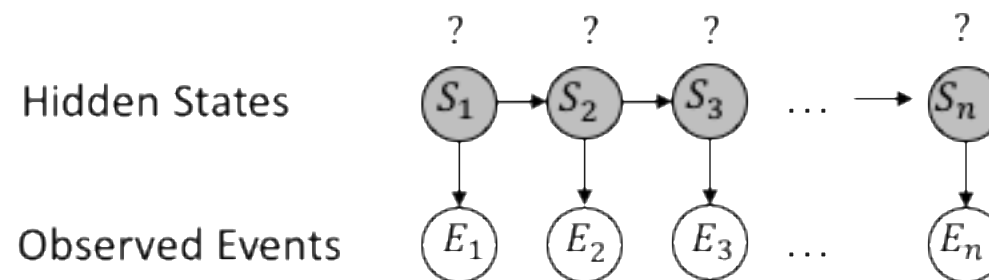
Bayes Network



Calculate the joint probability $P(B, J, A, E, M)$
 $= P(J|A) P(M|A) P(A|B, E) P(B) P(E)$

Calculate the state probability $P(B) = \sum_{J, A, E, M} P(B, J, A, E, M)$

Hidden Markov Model



Calculate the conditional distribution $P(S_t | E_1, \dots, E_n)$

Factorize (Bayes Theorem)

$$\propto P(S_t | E_1, \dots, E_t) * P(E_{t+1}, \dots, E_n | S_t, E_1, \dots, E_t)$$

Use the Markov Property

$$= P(S_t | E_1, \dots, E_t) * P(E_{t+1}, \dots, E_n | S_t)$$

Forward Backward Algorithm

$$= \alpha_t \odot \beta_t$$

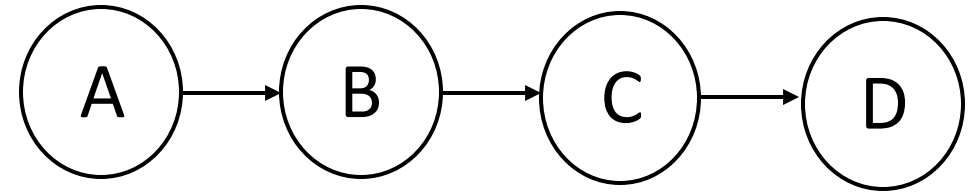
Just involves computation of joint distributions and its marginalization

Example of inference on a Bayesian Network

Consider the following Bayesian Network

- $A \in \{a^1, a^2\}, B \in \{b^1, b^2\}, C \in \{c^1, c^2\}, D \in \{d^1, d^2\}$

Inference task: Compute $P(D)$



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

- Simple way would be to generate each possible sequence (A,B,C,D) and sum over them
 - Exponential in the number of variables

Example of inference on Bayesian Network

- Enumerating all combinations
- Each term has 3 multiplications
- $8+8 = 16$ terms
- Total multiplication ops = $16 \times 3 = 48$
- 7 additions for d^1 and 7 additions for d^2
- Total additions ops = $7+7=14$

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

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

All terms involved in computation of $P(d^1)$ and $P(d^2)$ respectively.

Example of inference on Bayesian Network

Can we reduce the number of computations?

- Many terms are common; they can be computed once and reused

Consider the **orange highlighted** box, $P(c^1|b^1)P(d^1|c^1)$ is common.

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

Consider the blue **highlighted** box, $P(c^1|b^2)P(d^1|c^1)$ is common.

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

Define: $\tau_1(B) = P(a^1)P(B|a^1) + P(a^2)P(B|a^2)$

where $B \in \{b^1, b^2\}$

	$P(a^1)$	$P(b^1 a^1)$	$P(c^1 b^1)$	$P(d^1 c^1)$
+	$P(a^2)$	$P(b^1 a^2)$	$P(c^1 b^1)$	$P(d^1 c^1)$
+	$P(a^1)$	$P(b^2 a^1)$	$P(c^1 b^2)$	$P(d^1 c^1)$
+	$P(a^2)$	$P(b^2 a^2)$	$P(c^1 b^2)$	$P(d^1 c^1)$
+	$P(a^1)$	$P(b^1 a^1)$	$P(c^2 b^1)$	$P(d^1 c^2)$
+	$P(a^2)$	$P(b^1 a^2)$	$P(c^2 b^1)$	$P(d^1 c^2)$
+	$P(a^1)$	$P(b^2 a^1)$	$P(c^2 b^2)$	$P(d^1 c^2)$
+	$P(a^2)$	$P(b^2 a^2)$	$P(c^2 b^2)$	$P(d^1 c^2)$
	$P(a^1)$	$P(b^1 a^1)$	$P(c^1 b^1)$	$P(d^2 c^1)$
+	$P(a^2)$	$P(b^1 a^2)$	$P(c^1 b^1)$	$P(d^2 c^1)$
+	$P(a^1)$	$P(b^2 a^1)$	$P(c^1 b^2)$	$P(d^2 c^1)$
+	$P(a^2)$	$P(b^2 a^2)$	$P(c^1 b^2)$	$P(d^2 c^1)$
+	$P(a^1)$	$P(b^1 a^1)$	$P(c^2 b^1)$	$P(d^2 c^2)$
+	$P(a^2)$	$P(b^1 a^2)$	$P(c^2 b^1)$	$P(d^2 c^2)$
+	$P(a^1)$	$P(b^2 a^1)$	$P(c^2 b^2)$	$P(d^2 c^2)$
+	$P(a^2)$	$P(b^2 a^2)$	$P(c^2 b^2)$	$P(d^2 c^2)$

All terms involved in computation of $P(d^1)$ and $P(d^2)$ respectively.

Example of inference on Bayesian Network

Consider the **orange highlighted** box, $P(d^1|c^1)$ is common.

Compute: $\tau_1(b^1)P(c^1|b^1) + \tau_1(b^2)P(c^1|b^2)$

	$\tau_1(b^1)$	$P(c^1 b^1)$	$P(d^1 c^1)$
+	$\tau_1(b^2)$	$P(c^1 b^2)$	$P(d^1 c^1)$
+	$\tau_1(b^1)$	$P(c^2 b^1)$	$P(d^1 c^2)$
+	$\tau_1(b^2)$	$P(c^2 b^2)$	$P(d^1 c^2)$

Consider the **blue highlighted** box, $P(d^1|c^2)$ is common.

Compute: $\tau_1(b^1)P(c^2|b^1) + \tau_1(b^2)P(c^2|b^2)$

	$\tau_1(b^1)$	$P(c^1 b^1)$	$P(d^2 c^1)$
+	$\tau_1(b^2)$	$P(c^1 b^2)$	$P(d^2 c^1)$
+	$\tau_1(b^1)$	$P(c^2 b^1)$	$P(d^2 c^2)$
+	$\tau_1(b^2)$	$P(c^2 b^2)$	$P(d^2 c^2)$

Define: $\tau_2(C) = \tau_1(b^1)P(C|b^1) + \tau_1(b^2)P(C|b^2)$
 where $C \in \{c^1, c^2\}$

All terms involved in computation of $P(d^1)$ and $P(d^2)$ respectively. The sum is simplified because of use of $\tau_1(B)$.

Example of inference on Bayesian Network

- Computation shown alongside is easy and gives $P(D)$
- Previous steps are equivalent to **pushing summation inside**

$$P(D) = \sum_C \sum_B \sum_A P(A)P(B|A)P(C|B)P(D|C)$$

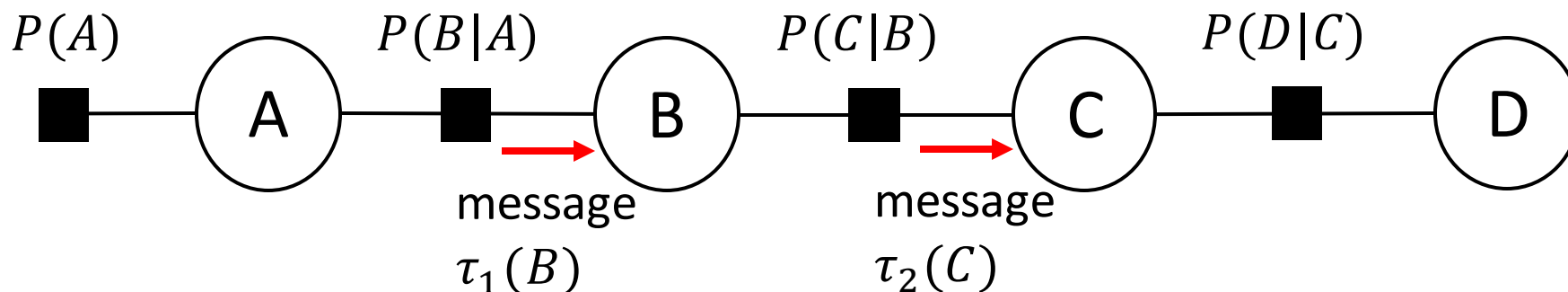
$$P(D) = \sum_C P(D|C) \sum_B P(C|B) \underbrace{\sum_A P(A)P(B|A)}_{\tau_1(B)} \underbrace{}_{\tau_2(C)}$$

$$+ \tau_2(c^1) P(d^1 | c^1) + \tau_2(c^2) P(d^1 | c^2)$$

$$+ \tau_2(c^1) P(d^2 | c^1) + \tau_2(c^2) P(d^2 | c^2)$$

Computation of $P(D)$ is simplified because of use of $\tau_1(B)$, $\tau_2(C)$.

Flow of computations (messages) in Factor Graph corresponding to the given BN



Sum-product algorithm reduces computations

- Pushing summations inside reduced the number of computations
 - Simple way: 48 multiplications + 14 additions
 - Pushing summations inside: 4x3 multiplications + 2x3 additions
 - Can be up to linear in number of variables (much better than exponential!)
- What helped in addressing the exponential blowup of marginalizing the joint distribution?
 - Graph structure – because of structure of Bayesian Network, some subexpression in the joint depend only on a small number of variables
 - Pushing summation inside – by computing these expressions once and caching the results, we can avoid generating them exponentially many times
- Referred to as **sum-product algorithm** or **Belief Propagation**

Inference problem on Factor Graphs

What is the problem we are trying to solve?

- Marginalization on Factor Graphs

Marginal probability

$$P(X_i) = \sum_{\mathbf{X} \setminus X_i} P(\mathbf{X})$$

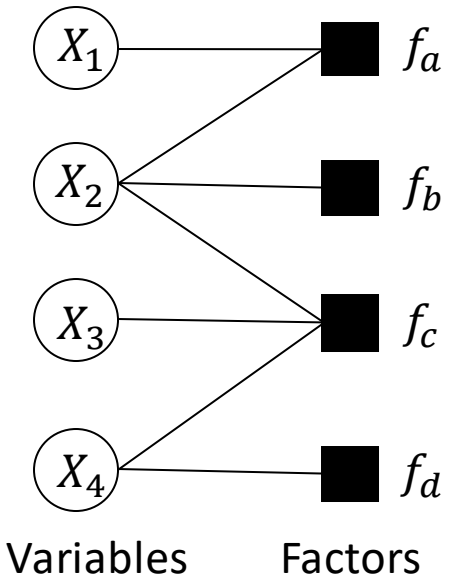
All variables except X_i

Challenge:

- Computationally expensive because the sum is calculated on all variables except one

Approach:

- Factorize the joint distribution according the structure
- Use belief propagation to reduce computations

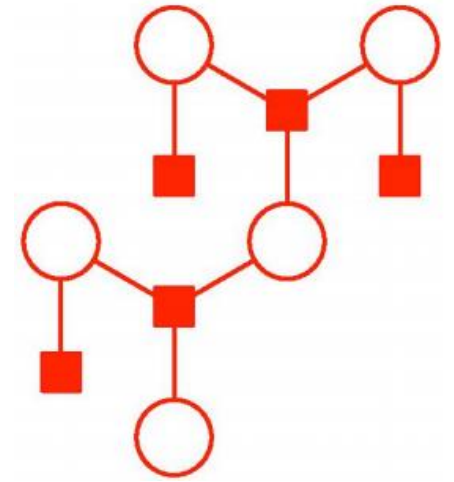


$X_i \in \{0,1\}$ is a discrete variable
e.g., a Boolean variable

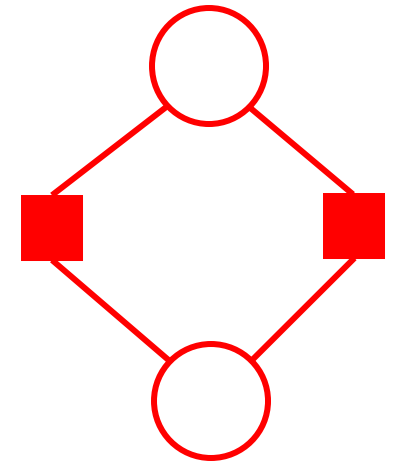
$f_c(X_c)$ is a tensor
on a set of variables X_c

Belief propagation

- Also known as sum-product algorithm
- Computes marginal distributions by “pushing in summations”
- Exact inference for linear graphs and trees
- Approximate inference for graphs with loops; performs remarkably well
- In case of Factor Graphs, involves two types of **messages**
 - From factor to variables
 - From variables to factors



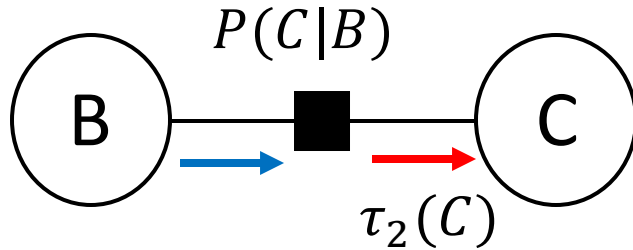
Tree Factor graph



Factor graph with loop

Belief Propagation – Message from factor to variable

Recall from previous example:



$$\tau_2(C) = \sum_B P(C|B) \tau_1(B)$$

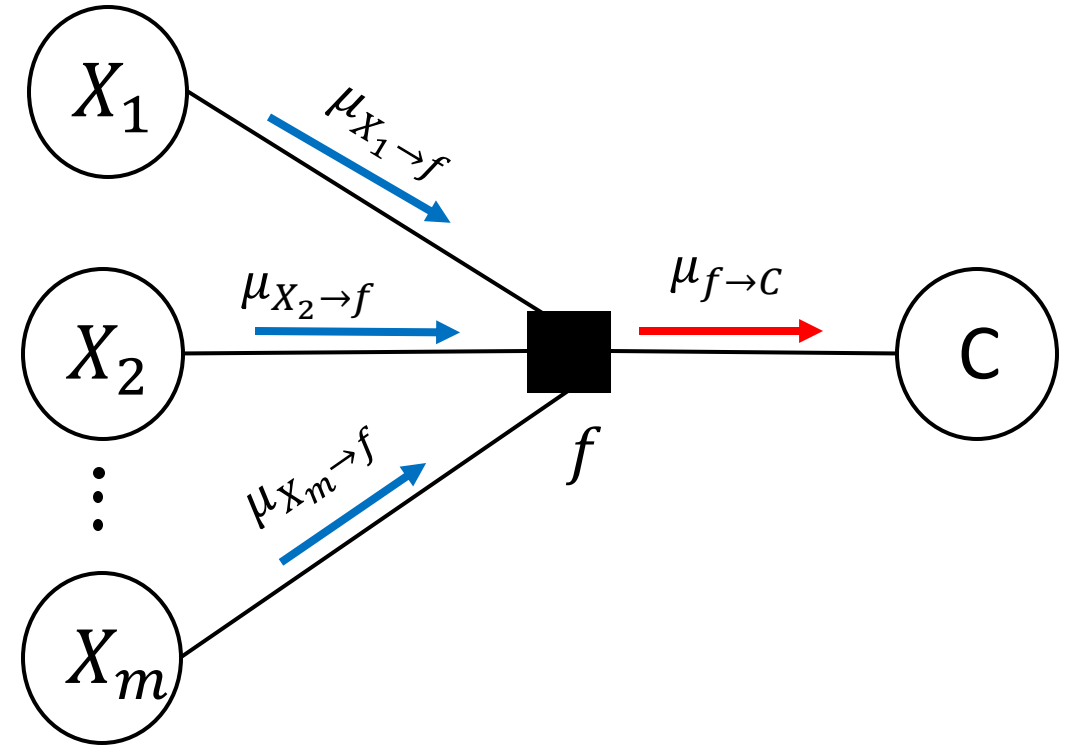
To get the general expression,
denote by:

$$f(B, C) = P(C|B)$$

$$\mu_{f \rightarrow C}(C) = \tau_2(C)$$

$$\mu_{B \rightarrow f}(B) = \tau_1(B)$$

In general:

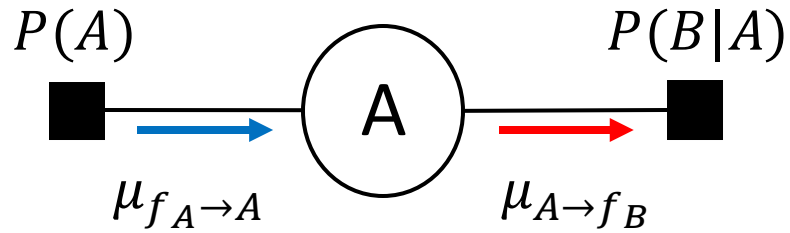


$$\mu_{f \rightarrow C}(C) = \sum_{X_1, X_2, \dots, X_m} f(C, X_1, \dots, X_m) \prod_{i=1}^m \mu_{X_i \rightarrow f}(X_i)$$

Message from factor to variable: Product of all incoming messages and factor, sum out previous variables

Belief Propagation – Message variable to factor

Recall from previous example:

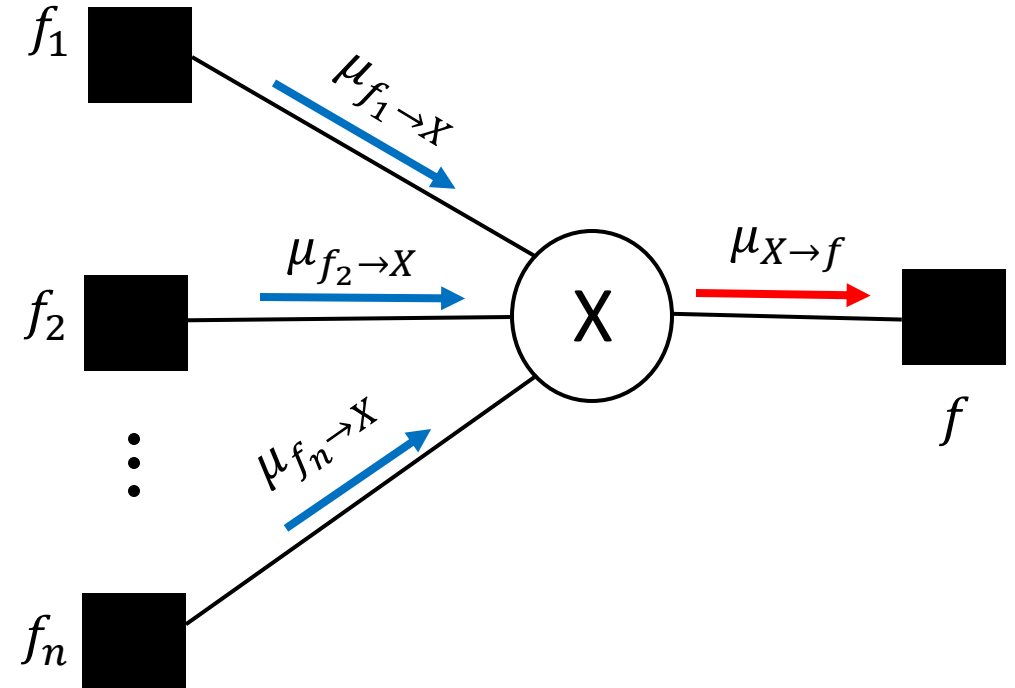


$$\mu_{A \rightarrow f_B}(A) = \mu_{f_A \rightarrow A}(A) = P(A)$$

Where,

$$\begin{aligned} f_A(A) &= P(A) \\ f_B(A, B) &= P(B|A) \end{aligned}$$

In general:



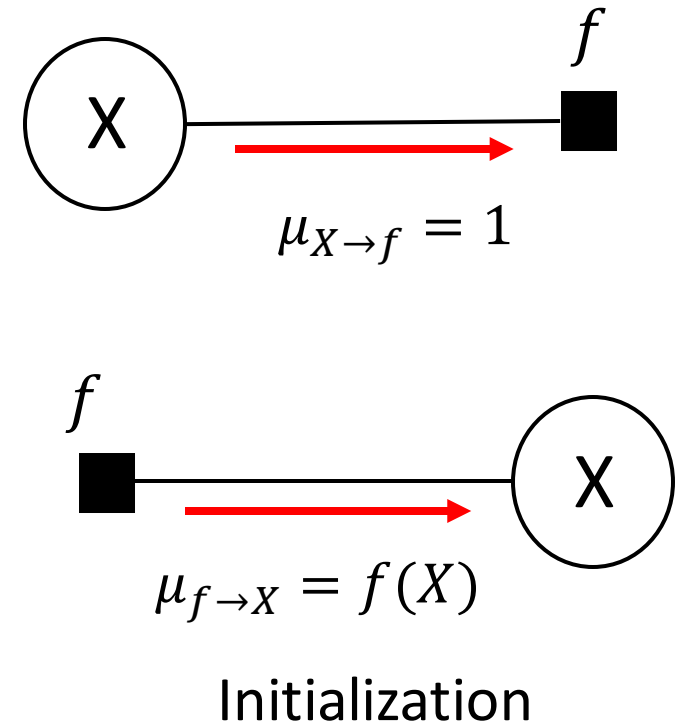
$$\mu_{X \rightarrow f}(X) = \prod_{i=1}^n \mu_{f_i \rightarrow X}(X)$$

Message from variable to factor: Product of all incoming messages

Belief Propagation: General Algorithm

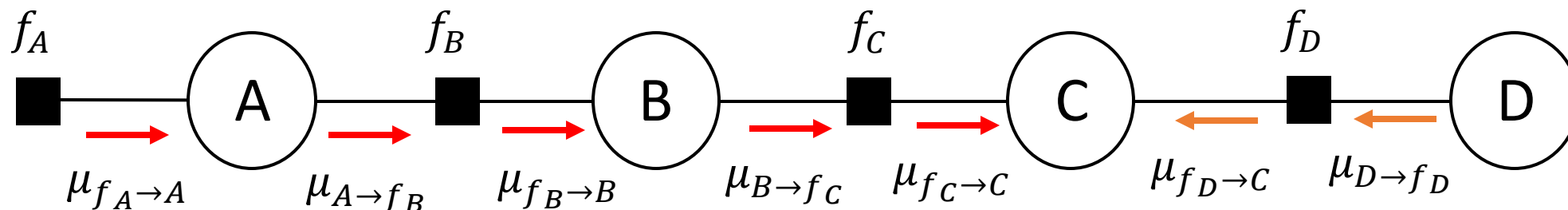
Steps to compute marginal distribution for all variables

- How to start the algorithm
 - Choose a node in the factor graph as root node
 - Compute all the leaf-to-root messages
 - Compute all the root-to-leaf messages
- Initial Conditions
 - Starting from a factor leaf/root node, the initial factor-to-variable message is the factor itself
 - Starting from a variable leaf/root node, the initial variable-to-factor message is a vector of ones
- Computing marginals
 - Marginal is given by the product of all incoming messages; normalize if necessary



Example of belief propagation

Compute $P(C)$



$$\mu_{f_A \rightarrow A}(A) = f_A(A) = P(A)$$

$$\mu_{A \rightarrow f_B}(A) = \mu_{f_A \rightarrow A}(A) = P(A)$$

$$\begin{aligned} \mu_{f_B \rightarrow B}(B) &= \sum_A f_B(A, B) \mu_{f_A \rightarrow A}(A) \\ &= \sum_A P(B|A) P(A) \end{aligned}$$

$$\mu_{B \rightarrow f_C}(B) = \mu_{f_B \rightarrow B}(B)$$

$$= \sum_A P(B|A) P(A)$$

$$\mu_{f_C \rightarrow C}(C) = \sum_B f_C(B, C) \mu_{B \rightarrow f_C}(B)$$

$$= \sum_B P(C|B) \sum_A P(B|A) P(A)$$

$$\mu_{D \rightarrow f_D}(D) = 1$$

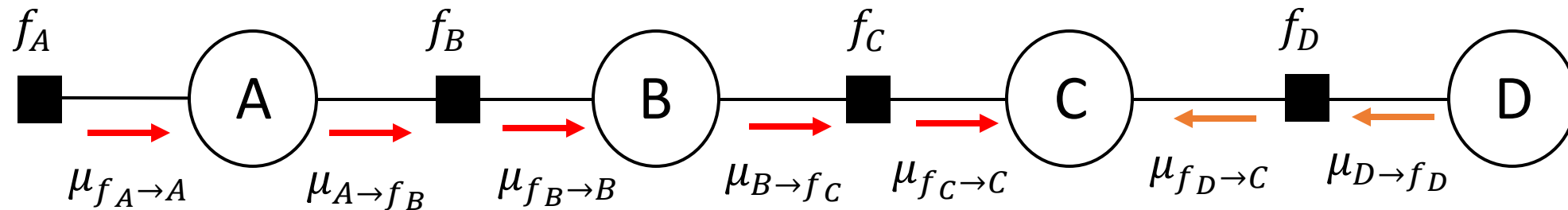
$$\mu_{f_D \rightarrow C}(C)$$

$$= \sum_D f_D(C, D) \mu_{f_D \rightarrow C}(C)$$

$$= \sum_D P(D|C)$$

Example of belief propagation

Compute $P(C)$



$$\mu_{f_D \rightarrow C}(C) = \sum_D P(D|C)$$

$$\mu_{f_C \rightarrow C}(C) = \sum_B P(C|B) \sum_A P(B|A)P(A)$$

$$P(C) = \mu_{f_C \rightarrow C}(C) \mu_{f_D \rightarrow C}(C)$$

Verifying that the above computation gives the marginal distribution

$$\begin{aligned} P(C) &= \left(\sum_B P(C|B) \sum_A P(B|A)P(A) \right) \left(\sum_D P(D|C) \right) = \sum_B P(C|B) \sum_A P(B|A)P(A) \sum_D P(D|C) \\ &= \sum_A \sum_B \sum_D P(A)P(B|A)P(C|B)P(D|C) = \sum_{A,B,D} P(A, B, C, D) \end{aligned}$$

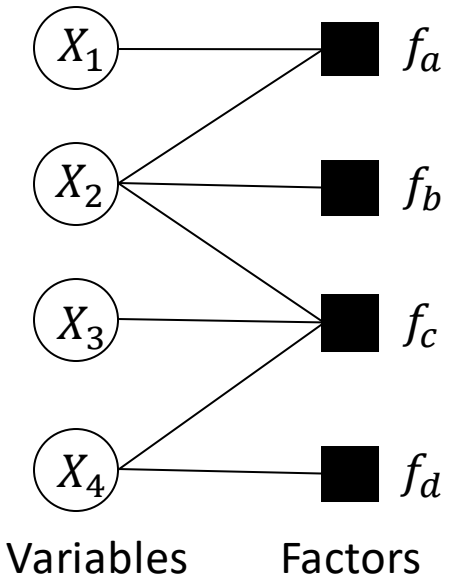
Marginalization on Factor Graphs

Marginal probability

$$P(X_i) = \sum_{\mathbf{X} \setminus X_i} P(\mathbf{X})$$

All variables except X_i

Inference method	Description
Belief Propagation	Exact inference on non-loop FG
Sampling - Markov Chain Monte Carlo, Gibbs	Approximate inference
Variational Inference	Approximate inference



$X_i \in \{0,1\}$ is a discrete variable
e.g., a Boolean variable

$f_c(X_c)$ is a tensor
on a set of variables X_c

References

- Daphne Koller, Nir Friedman's textbook on Graphical Models
- https://www.doc.ic.ac.uk/~mpd37/teaching/ml_tutorials/2016-11-09-Svensson-BP.pdf