

Unsupervised Learning

Graphical Models

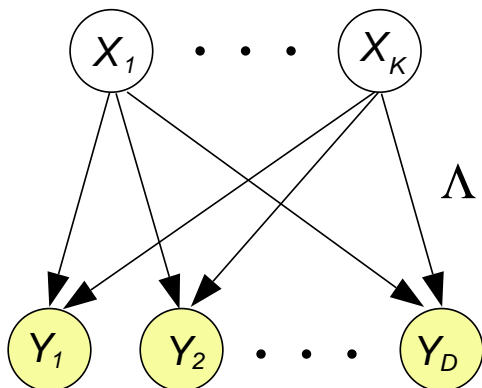
Zoubin Ghahramani

`zoubin@gatsby.ucl.ac.uk`

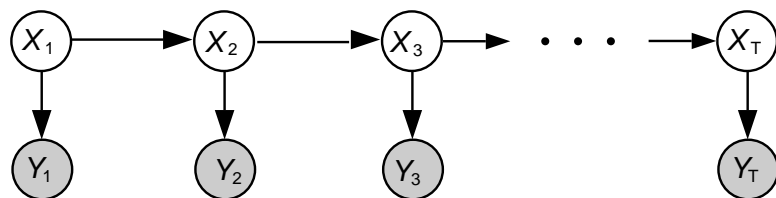
**Gatsby Computational Neuroscience Unit, and
MSc in Intelligent Systems, Dept Computer Science
University College London**

Autumn 2003

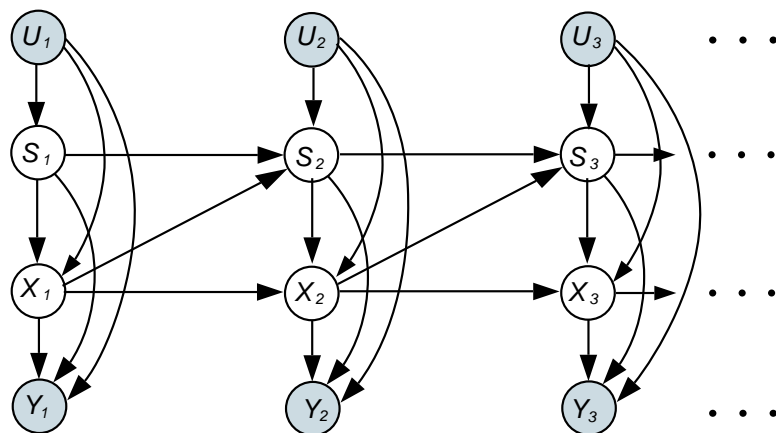
Some Examples



factor analysis
probabilistic PCA
ICA

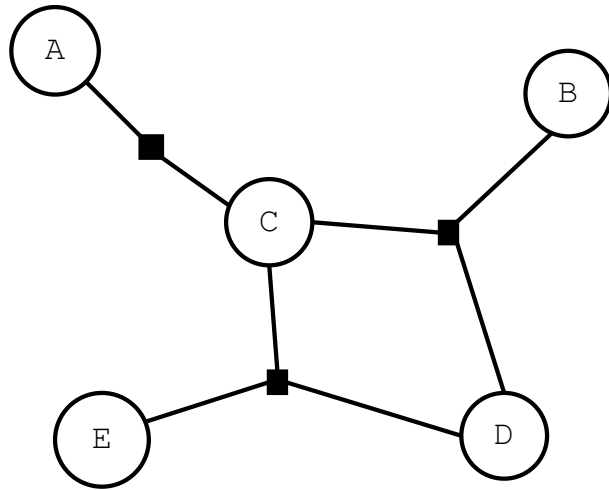


hidden Markov models
linear dynamical systems

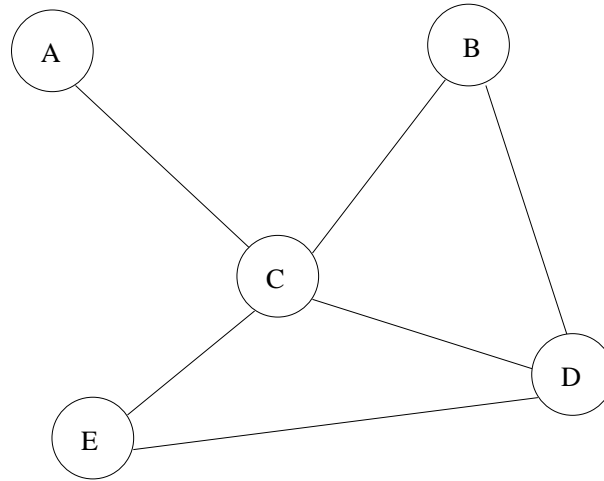


switching state-space models

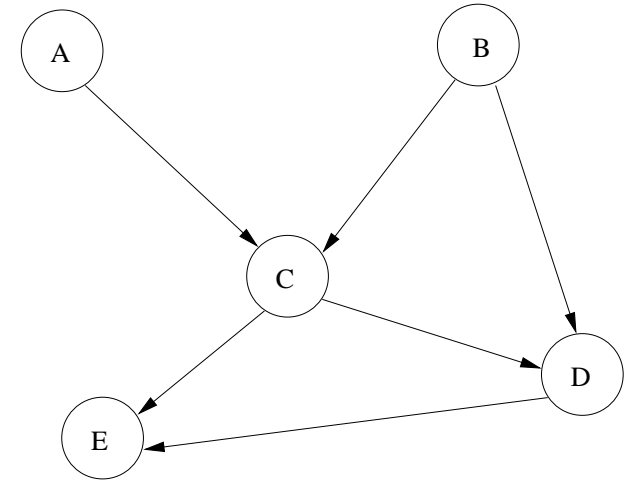
Three kinds of graphical models



factor graph



undirected graph



directed graph

Why do we need graphical models?

- Graphs are an **intuitive** way of representing and visualising the relationships between many variables. (Examples: family trees, electric circuit diagrams, neural networks)
- A graph allows us to abstract out the **conditional independence** relationships between the variables from the details of their parametric forms. Thus we can ask questions like: “Is A dependent on B given that we know the value of C ?” just by looking at the graph.
- Graphical models allow us to define general **message-passing algorithms** that implement Bayesian inference efficiently. Thus we can answer queries like “What is $P(A|C = c)$?” without enumerating all settings of all variables in the model.

Conditional Independence

Conditional Independence:

$$X \perp\!\!\!\perp Y|V \Leftrightarrow p(X|Y, V) = p(X|V)$$

when $p(Y, V) > 0$. Also

$$X \perp\!\!\!\perp Y|V \Leftrightarrow p(X, Y|V) = p(X|V)p(Y|V)$$

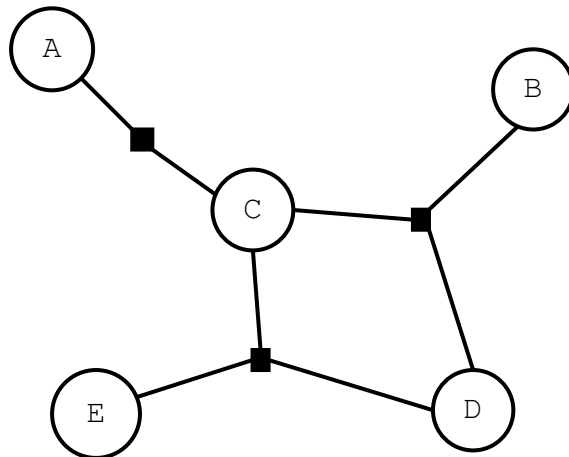
In general we can think of conditional independence between **sets of variables**:

$$\mathcal{X} \perp\!\!\!\perp \mathcal{Y}|\mathcal{V} \Leftrightarrow \{X \perp\!\!\!\perp Y|\mathcal{V}, \forall X \in \mathcal{X} \text{ and } \forall Y \in \mathcal{Y}\}$$

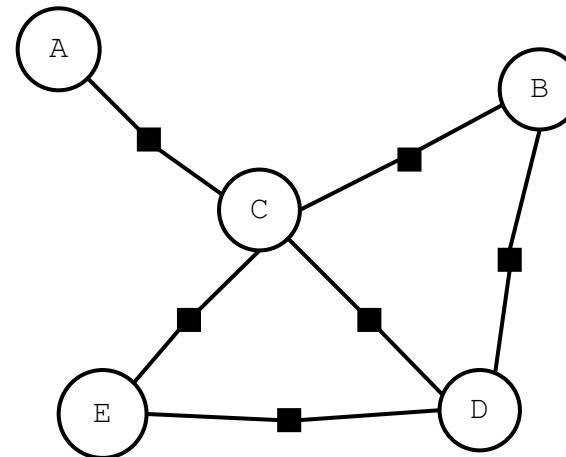
Marginal Independence:

$$X \perp\!\!\!\perp Y \Leftrightarrow X \perp\!\!\!\perp Y|\emptyset \Leftrightarrow p(X, Y) = p(X)p(Y)$$

Factor Graphs



(a)



(b)

The circles in a factor graph represent random variables.
The filled dots represent factors in the joint distribution.

$$(a) P(A, B, C, D, E) = \frac{1}{Z} g_1(A, C) g_2(B, C, D) g_3(C, D, E)$$

$$(b) P(A, B, C, D, E) = \frac{1}{Z} g_1(A, C) g_2(B, C) g_3(C, D) g_4(B, D) g_5(C, E) g_6(D, E)$$

The g_i are non-negative functions of their arguments, and Z is a normalization constant.
Two nodes are **neighbors** if they share a common factor.

Fact: $X \perp\!\!\!\perp Y \mid \mathcal{V}$ if every path between X and Y contains some node $V \in \mathcal{V}$

Corollary: Given the neighbors of X , the variable X is **conditionally independent** of all other variables: $X \perp\!\!\!\perp Y \mid \text{ne}(X)$, $\forall Y \notin \{X \cup \text{ne}(X)\}$

What is an Undirected Graphical Model?

In an Undirected Graphical Model (or Markov Network), the joint probability over all variables can be written in a factored form:

$$P(\mathbf{x}) = \frac{1}{Z} \prod_j g_j(\mathbf{x}_{C_j})$$

where $\mathbf{x} = [x_1, \dots, x_K]$, and

$$C_j \subseteq \{1, \dots, K\}$$

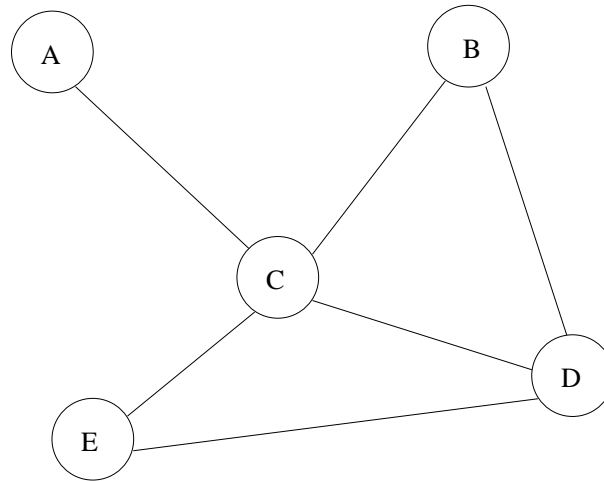
are subsets of the set of all variables, and $\mathbf{x}_S \equiv [x_k : k \in S]$.

This type of probabilistic model can be represented **graphically**.

Graph Definition: Let each variable be a node. Connect nodes i and k if there exists a set C_j such that both $i \in C_j$ and $k \in C_j$. These sets form the *cliques* of the graph (fully connected subgraphs).

Note: Undirected Graphical Models are also called *Markov Networks*.

Undirected Graphical Models (Markov Networks)



$$P(A, B, C, D, E) = \frac{1}{Z} g_1(A, C) g_2(B, C, D) g_3(C, D, E)$$

Fact: $X \perp\!\!\!\perp Y \mid \mathcal{V}$ if every path between X and Y contains some node $V \in \mathcal{V}$

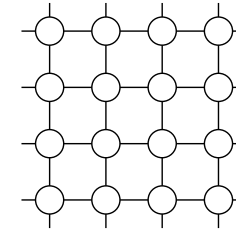
Corollary: Given the neighbors of X , the variable X is conditionally independent of all other variables: $X \perp\!\!\!\perp Y \mid \text{ne}(X)$, $\forall Y \notin \{X \cup \text{ne}(X)\}$

Markov Blanket: \mathcal{V} is a Markov Blanket for X iff $X \perp\!\!\!\perp Y \mid \mathcal{V}$ for all $Y \notin \{X \cup \mathcal{V}\}$.

Markov Boundary: minimal Markov Blanket $\equiv \text{ne}(X)$ for undirected graphs

Examples of Undirected Graphical Models

- Markov Random Fields (used in Computer Vision)



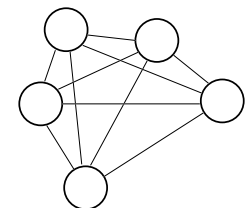
- Exponential Language Models (used in Speech and Language Modelling)

$$p(s) = \frac{1}{Z} p_0(s) \exp \left\{ \sum_i \lambda_i f_i(s) \right\}$$

- Products of Experts (widely applicable)

$$p(\mathbf{x}) = \frac{1}{Z} \prod_j p_j(\mathbf{x}|\theta_j)$$

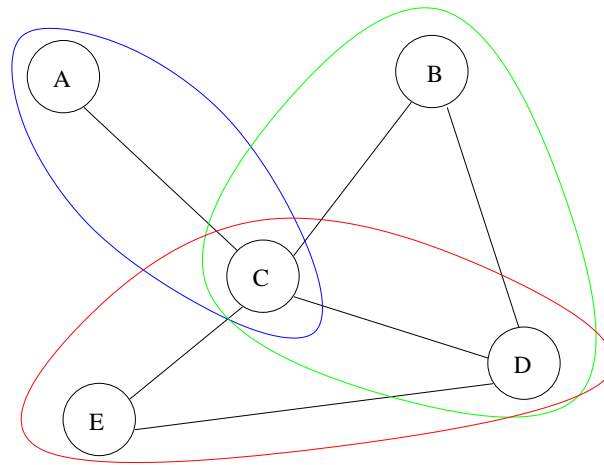
- Boltzmann Machines (a kind of Neural Network/Ising Model)



Clique Potentials and Undirected Graphs (Markov Networks)

Definition: a *clique* is a fully connected subgraph. By clique we usually mean maximal clique (i.e. not contained within another clique)

C_i will denote the set of variables in the i^{th} clique.



1. Identify cliques of graph G
2. For each clique C_i assign a non-negative function $g_i(\mathbf{x}_{C_i})$ which measures “compatibility”.

3. $p(x_1, \dots, x_K) = \frac{1}{Z} \prod_i g_i(\mathbf{x}_{C_i})$ where $Z = \sum_{x_1 \dots x_K} \prod_i g_i(\mathbf{x}_{C_i})$ is the normalization

If V lies in *all* paths between X and Y in G , then $X \perp\!\!\!\perp Y | V$ in p .

Hammersley–Clifford Theorem (1971)

Theorem: A probability function p formed by a normalized product of positive functions on cliques of G is a Markov Field relative to G .

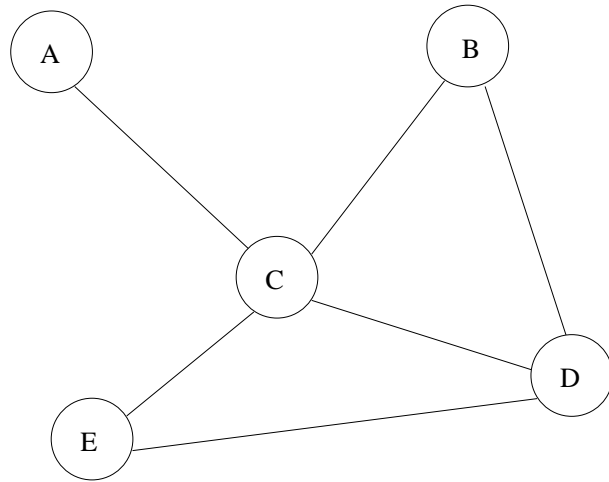
Definition: The distribution p is a *Markov Field relative to G* if G does not imply any conditional independence relationships that are not true in p .
(We are usually interested in the minimal such graph.)

Proof: We need to show that if p is a product of functions on cliques of G then a variable is conditionally independent of its non-neighbors in G given its neighbors in G . That is: $\text{ne}(x_\ell)$ is a Markov Blanket for x_ℓ :

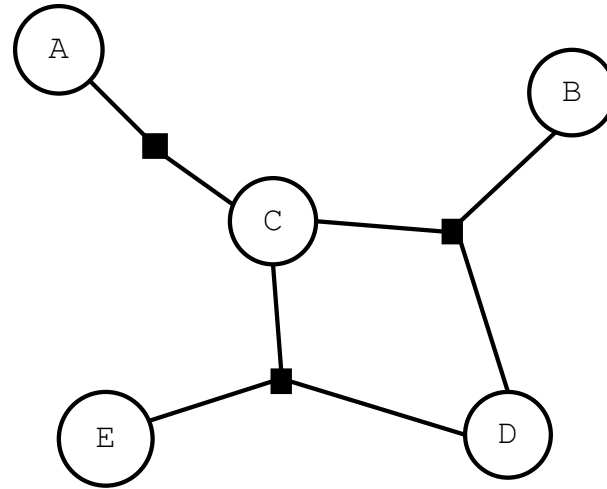
$$\begin{aligned} p(x_\ell, x_m, \dots) &= \frac{1}{Z} \prod_i g_i(\mathbf{x}_{C_i}) = \frac{1}{Z} \prod_{i:\ell \in C_i} g_i(\mathbf{x}_{C_i}) \prod_{j:\ell \notin C_j} g_j(\mathbf{x}_{C_j}) \\ &= \frac{1}{Z} f_1(x_\ell, \text{ne}(x_\ell)) f_2(\text{ne}(x_\ell), x_m) = \frac{1}{Z'} p(x_\ell | \text{ne}(x_\ell)) p(x_m | \text{ne}(x_\ell)) \end{aligned}$$

This shows that: $p(x_\ell, x_m | \text{ne}(x_\ell)) = p(x_\ell | \text{ne}(x_\ell)) p(x_m | \text{ne}(x_\ell)) \Leftrightarrow x_\ell \perp\!\!\!\perp x_m | \text{ne}(x_\ell)$.

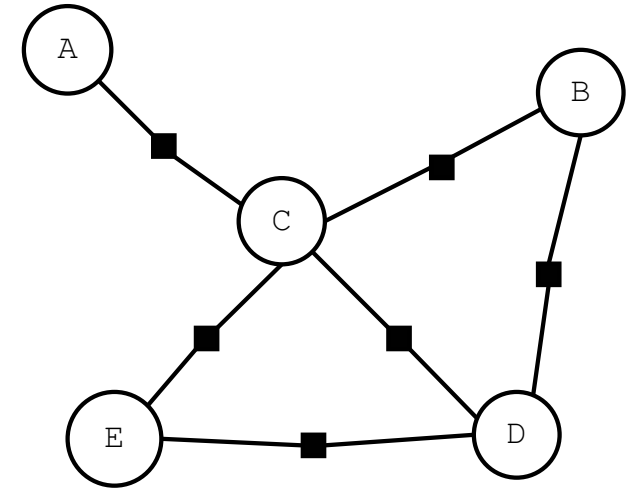
Comparing Undirected Graphs and Factor Graphs



(a)



(b)



(c)

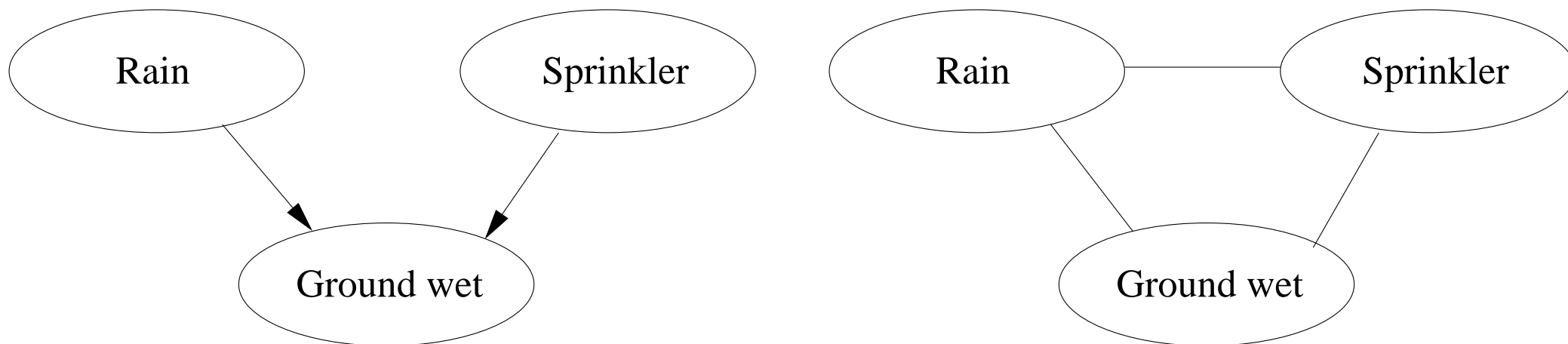
All nodes in (a), (b), and (c) have exactly the same neighbors and therefore these three graphs represent exactly the same conditional independence relationships.

(c) also represents the fact that the probability factors into a product of pairwise functions.

Consider the case where each variables is discrete and can take on K possible values. Then the functions in (a) and (b) are tables with $\mathcal{O}(K^3)$ cells, whereas in (c) they are $\mathcal{O}(K^2)$.

Problems with Undirected Graphs and Factor Graphs

Many useful independencies are unrepresented — two variables are connected merely because some other variable depends on them:

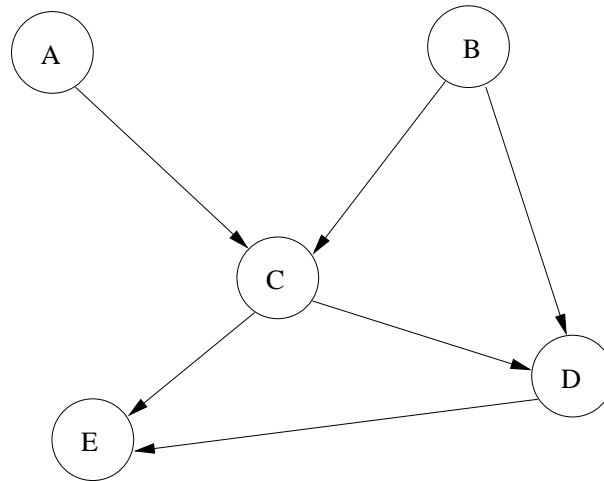


This highlights the difference between **marginal independence** and **conditional independence**.

R and S are marginally independent (i.e. given nothing), but they are conditionally dependent given G

“Explaining Away”: Observing that the spinkler is on, explains away the fact that the ground was wet, therefore we don’t need to believe that it rained.

Directed Acyclic Graphical Models (Bayesian Networks)



A DAG Model / Bayesian network corresponds to a factorization of the joint probability distribution:

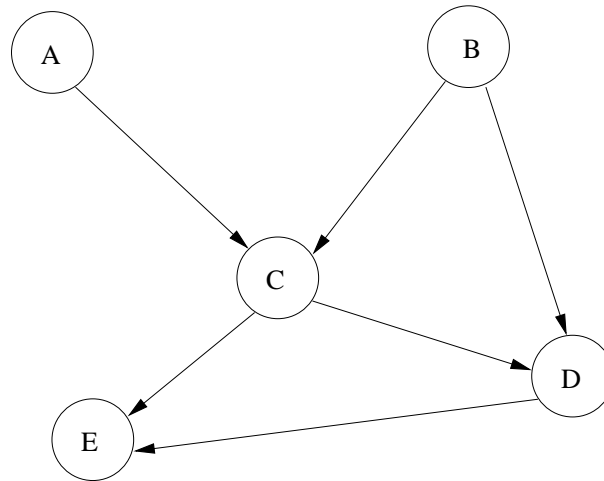
$$p(A, B, C, D, E) = p(A)p(B)p(C|A, B)p(D|B, C)p(E|C, D)$$

In general:

$$p(X_1, \dots, X_n) = \prod_{i=1}^n p(X_i | X_{\text{pa}(i)})$$

where $\text{pa}(i)$ are the parents of node i .

Directed Acyclic Graphical Models (Bayesian Networks)



Semantics: $X \perp\!\!\!\perp Y | \mathcal{V}$ if \mathcal{V} **d-separates** X from Y .

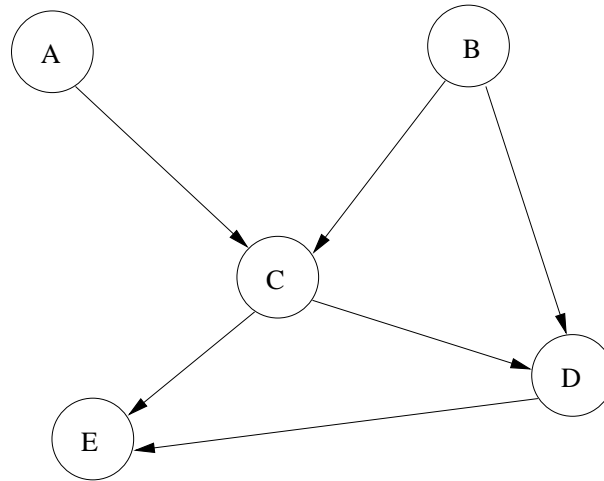
Definition: \mathcal{V} **d-separates** X from Y if every undirected path between X and Y is **blocked** by \mathcal{V} . A path is blocked if there is a node W on the path such that either:

1. W has converging arrows along the path ($\rightarrow W \leftarrow$) and neither W nor its descendants are in \mathcal{V} , or
2. W does not have converging arrows along the path ($\rightarrow W \rightarrow$ or $\leftarrow W \rightarrow$) and $W \in \mathcal{V}$.

Note that converging arrows *along the path* only refers to what happens on that path.

Corollary: Markov Boundary for X : $\{\text{parents}(X) \cup \text{children}(X) \cup \text{parents-of-children}(X)\}$.

The “Bayes-ball” algorithm

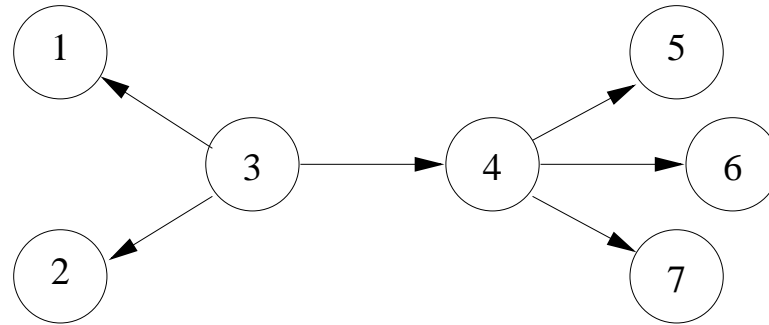


Game: can you get a ball from X to Y without being blocked by \mathcal{V} ?

Depending on the direction the ball came from and the type of node, the ball can **pass through** (from a parent to all children, from a child to all parents), **bounce back** (from any parent to all parents, or from any child to all children), or be **blocked**.

- An unobserved (hidden) node ($W \notin \mathcal{V}$) passes balls through but also bounces back balls from children.
- An observed (given) node ($W \in \mathcal{V}$) bounces back balls from parents but blocks balls from children.

From Directed Trees to Undirected Trees



$$p(1, 2, \dots, 7) = p(3)p(1|3)p(2|3)p(4|3)p(5|4)p(6|4)p(7|4)$$

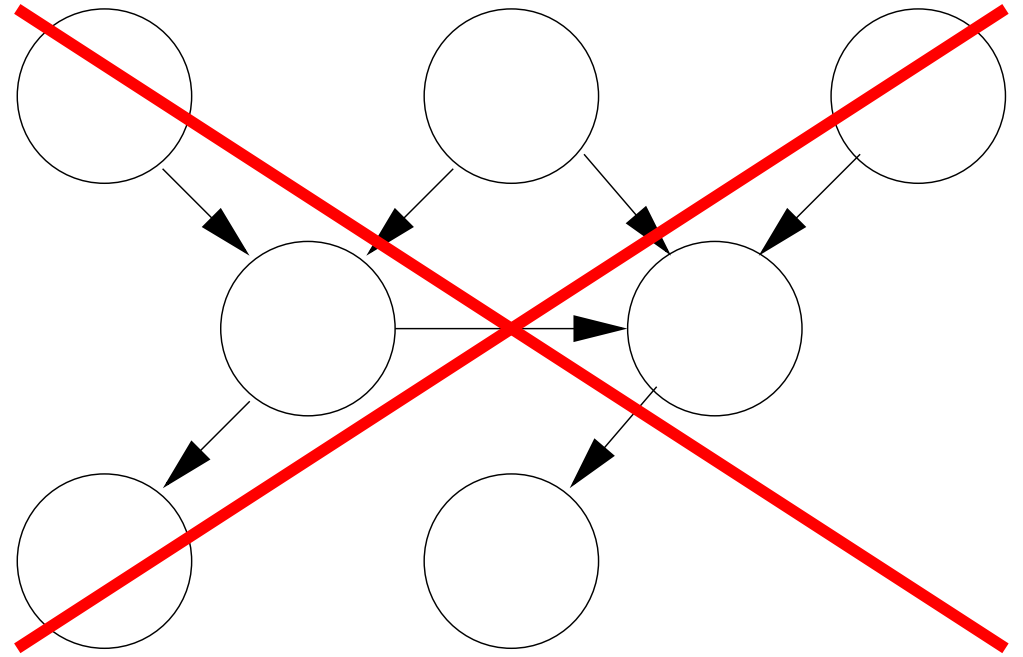
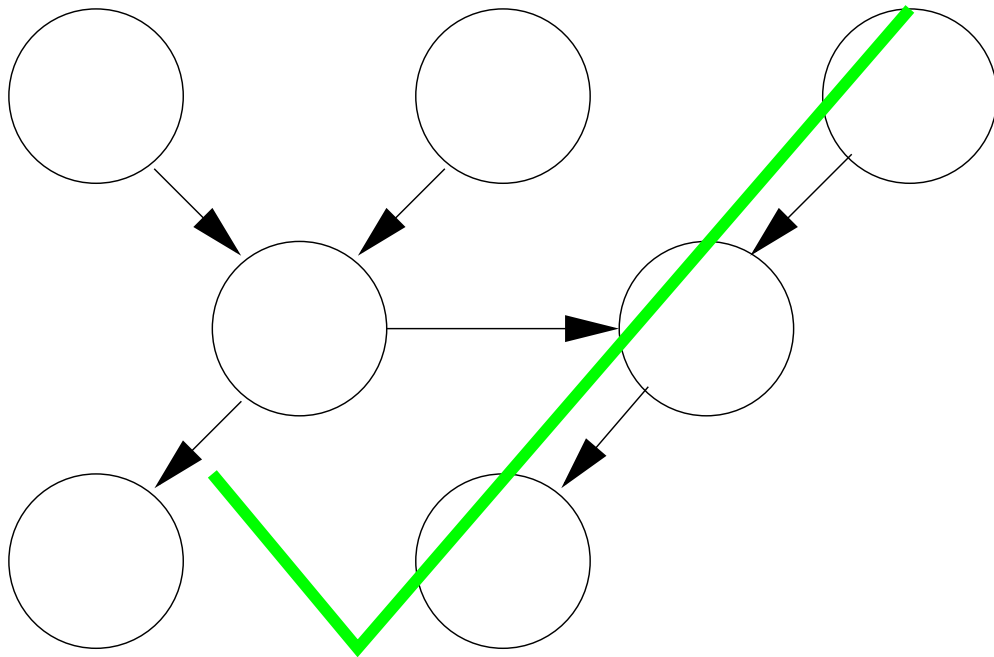
$$= \frac{p(1, 3)p(2, 3)p(3, 4)p(4, 5)p(4, 6)p(4, 7)}{p(3)p(3)p(4)p(4)p(4)}$$

$$= \frac{\text{product of cliques}}{\text{product of clique intersections}}$$

$$= g(1, 3)g(2, 3)g(3, 4)g(4, 5)g(4, 6)g(4, 7) = \prod_i g_i(C_i)$$

Belief Propagation (in Singly Connected Bayesian Networks)

Definition: S.C.B.N. has an undirected underlying graph which is a tree, *ie* there is only one path between any two nodes.



Goal: For some node X we want to compute $p(X|e)$ given evidence e .
Since we are considering S.C.B.N.s:

- every node X divides the evidence into **upstream** e_X^+ and **downstream** e_X^-
- every edge $X \rightarrow Y$ divides the evidence into **upstream** e_{XY}^+ and **downstream** e_{XY}^- .

The three key ideas behind Belief Propagation

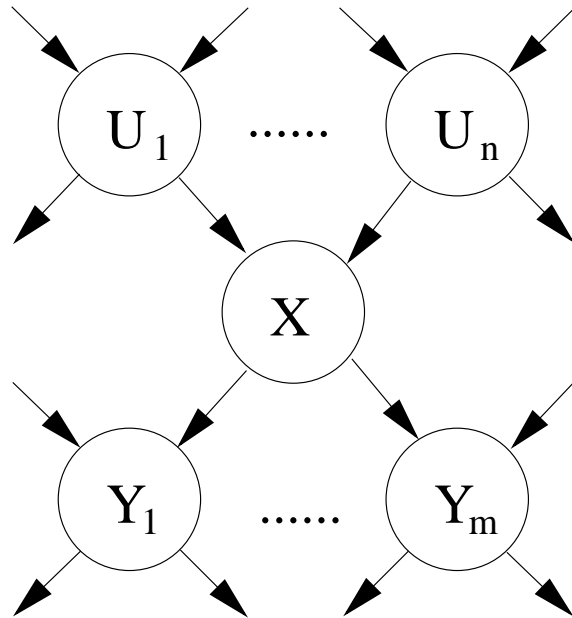
Idea 1: Our belief about the variable X can be found by combining upstream and downstream evidence:

$$\begin{aligned} p(X|e) &= \frac{p(X, e)}{p(e)} = \frac{p(X, e_X^+, e_X^-)}{p(e_X^+, e_X^-)} \propto p(X|e_X^+) \times \underbrace{p(e_X^-|X, e_X^+)}_{X \text{ d-separates } e_X^- \text{ from } e_X^+} \\ &= p(X|e_X^+)p(e_X^-|X) = \pi(X)\lambda(X) \end{aligned}$$

Idea 2: The upstream and downstream evidence can be computed via a local message passing algorithm between the nodes in the graph.

Idea 3: “Don’t send back to a node (any part of) the message it sent to you!”

Belief Propagation



top-down causal support:

$$\pi_X(U_i) = p(U_i | e_{U_i X}^+)$$

bottom-up diagnostic support:

$$\lambda_{Y_j}(X) = p(e_{XY_j}^- | X)$$

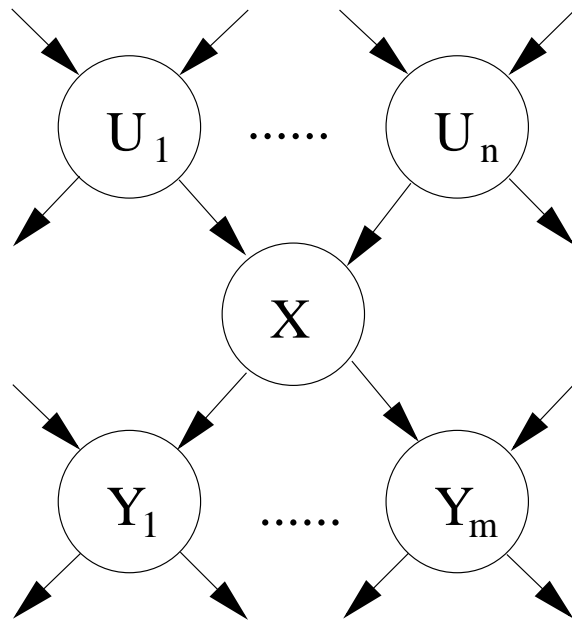
To update the belief about X :

$$\text{BEL}(X) = \frac{1}{Z} \lambda(X) \pi(X)$$

$$\lambda(X) = \prod_j \lambda_{Y_j}(X)$$

$$\pi(X) = \sum_{U_1 \dots U_n} p(X | U_1, \dots, U_n) \prod_i \pi_X(U_i)$$

Belief Propagation (cont.)



top-down causal support:

$$\pi_X(U_i) = p(U_i | e_{U_i X}^+)$$

bottom-up diagnostic support:

$$\lambda_{Y_j}(X) = p(e_{XY_j}^- | X)$$

Bottom-up propagation, message X sends to U_i :

$$\lambda_X(U_i) = \sum_X \lambda(X) \sum_{U_k: k \neq i} p(X | U_1, \dots, U_n) \prod_{k \neq i} \pi_X(U_k)$$

Top-down propagation, message X sends to Y_j :

$$\pi_{Y_j}(X) = \frac{1}{Z} \left[\prod_{k \neq j} \lambda_{Y_k}(X) \right] \sum_{U_1 \dots U_n} p(X | U_1, \dots, U_n) \prod_i \pi_X(U_i) = \frac{1}{Z} \frac{\text{BEL}(X)}{\lambda_{Y_j}(X)}$$

Z is the normaliser ensuring $\sum_X \pi_{Y_j}(X) = 1$

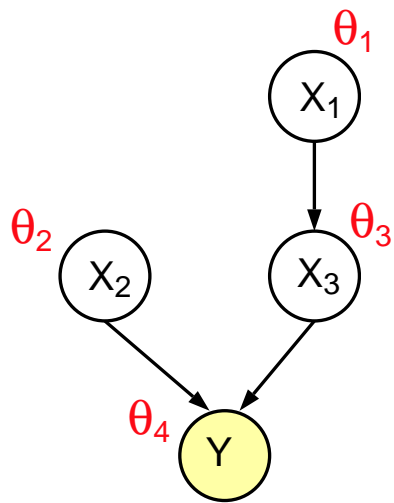
Belief Propagation in multiply connected Bayesian Networks

The Junction Tree algorithm: Form an undirected graph from your directed graph such that no additional conditional independence relationships have been created (this step is called “moralization”). Lump variables in cliques together and form a tree of cliques—this may require a nasty step called “triangulation”. Do inference in this tree.

Cutset Conditioning: or “reasoning by assumptions”. Find a small set of variables which, if they were given (i.e. known) would render the remaining graph singly connected. For each value of these variables run belief propagation on the singly connected network. Average the resulting beliefs with the appropriate weights.

Loopy Belief Propagation: just use BP although there are loops. In this case the terms “upstream” and “downstream” are not clearly defined. No guarantee of convergence, but often works well in practice.

Learning with Hidden Variables: The EM Algorithm



Assume a model parameterised by θ with observable variables Y and hidden variables X

Goal: maximise parameter log likelihood given observables.

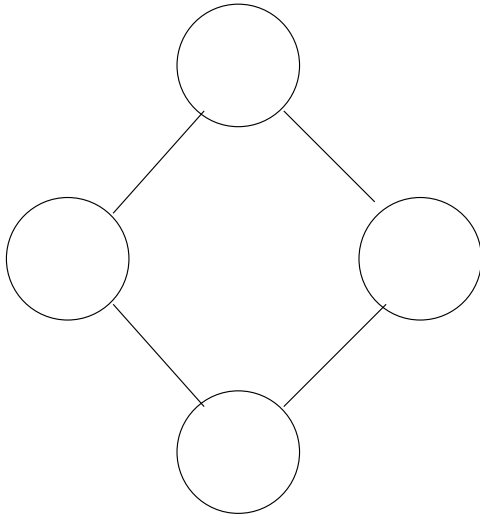
$$\mathcal{L}(\theta) = \ln p(Y|\theta) = \ln \sum_X p(Y, X|\theta)$$

- **E-step:** first infer $p(X|Y, \theta_{old})$, then
- **M-step:** find θ_{new} using complete data learning

The E-step requires solving the *inference* problem: finding explanations, X , for the data, Y , given the current model, θ (using e.g. BP).

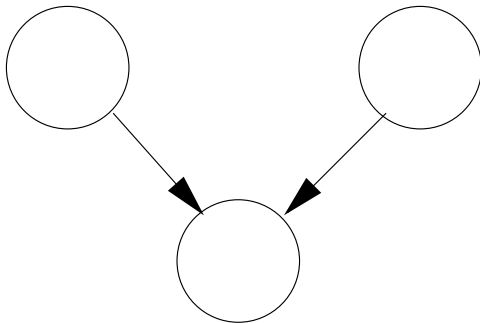
How about structure learning?

Expressive Power of Directed and Undirected Graphs



No Directed Graph (Bayesian network) can represent these and only these independencies

No matter how we direct the arrows there will always be two non-adjacent parents sharing a common child \implies dependence in Directed Graph but independence in Undirected Graph.



No Undirected Graph or Factor Graph can represent these and only these independencies