

Graphical Models

Seung-Hoon Na

Chonbuk National University

Directed Graphical Model

- **Belief networks, Bayesian belief networks**
- Convenient framework for representing independence assumptions

$$p(x_1, \dots, x_D) = \prod_{i=1}^D p(x_i | \text{pa}(x_i))$$

$\text{pa}(x_i)$ represent the parental variables of variable x_i

The Challenge of Unstructured Modeling

- With probabilistic models, we can do many tasks like
 - Density estimation
 - Denoising
 - Missing value imputation
 - Sampling
- These tasks are often more complicated than classification
- Learning probabilistic models requires a complete understanding of the entire structure of the input

Sampling



The Challenge of Unstructured Modeling

- Modeling a rich distribution over thousands or millions of random variables is a challenging task, both computationally and statistically
 - 32 × 32 pixel color (RGB) image
 - 2^{3072} binary images of this form → 10^{800} times larger than the estimated number of atoms in the universe
- Table-based approach
 - Representing $P(\mathbf{x})$ by storing a lookup table with one probability value per possible outcome
 - Requires k^n parameters!

Table-Based Model

- Specifying $p(x_1; \dots; x_N)$ over binary variables x_i
 - Takes $O(2^N)$ space \rightarrow impractical
 - computing a marginal - $p(x_1)$: summing over the 2^{N-1} states of the other variables

Table-Based Model

- Table-based model is not feasible because of
 - **Memory: the cost of storing the representation**
 - **Statistical efficiency**
 - Require an astronomically large training set to fit accurately, given an astronomical number of parameters
 - **Runtime: the cost of inference**
 - Computing the marginal distribution $P(x_1)$ or the conditional distribution $P(x_2 | x_1)$ require summing across the entire table
 - **Runtime: the cost of sampling**
 - Sample some value $u \sim U(0, 1)$, iterate through the table adding up the probability values until they exceed u
 - Requires reading through the whole table in the worst case

Table-based model

- Explicitly model every possible kind of interaction between every possible subset of variables
- The probability distributions we encounter in real tasks are much simpler than this.
- Usually, most variables influence each other only indirectly

Relay race example

- Modeling the finishing times of a team in a relay race
 - Three runners: Alice, Bob and Carol.
 - Model each of their finishing times as a continuous random variable
- Carol's finishing time depends only *indirectly* on Alice's finishing time via Bob's
 - ➔ We can model the relay race using only two interactions: Alice's effect on Bob and Bob's effect on Carol.

Structured Probabilistic Model

- Structure
 - constrain the nature of the variable interactions
 - specify which variables are independent of others, leading to a structured factorisation of the joint probability distribution.
 - E.g.) $p(x_1; \dots; x_{100}) = \prod_i \phi(x_i, x_{i+1})$

Structured probabilistic models (graphical models)

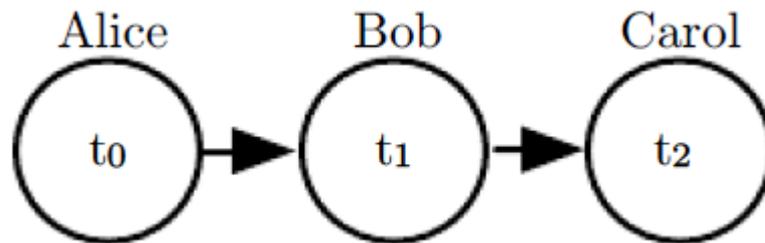
- Provide a **formal framework for modeling only direct interactions between random variables**
 - Allows the models to have significantly fewer parameters which can in turn be estimated reliably from less data
 - Reducing computational cost in terms of storing the model, performing inference in the model, and drawing samples from the model
- Graphical models can be divided into two categories:
 - 1) **Directed Graphical Models**
 - Based on directed acyclic graph
 - 2) **Undirected Graphical Models**
 - Based on undirected graphs

Directed Graphical Models

- known as the *belief network* or *Bayesian network*
- Defined on variables \mathbf{x} is defined by a directed acyclic graph G
 - With a set of *local conditional probability distributions*
- The probability distribution over \mathbf{x} :

$$p(\mathbf{x}) = \prod_i p(x_i \mid Pa_G(x_i))$$

Relay Race Example



$$p(t_0, t_1, t_2) = p(t_0)p(t_1 | t_0)p(t_2 | t_1)$$

- Discretizing time: t_0 , t_1 and t_2 each be discrete variables with 100 possible values
- Unstructured model: $p(t_0, t_1, t_2)$ requires 999,999 values
- Structured model: requires a total of 19,899 values.

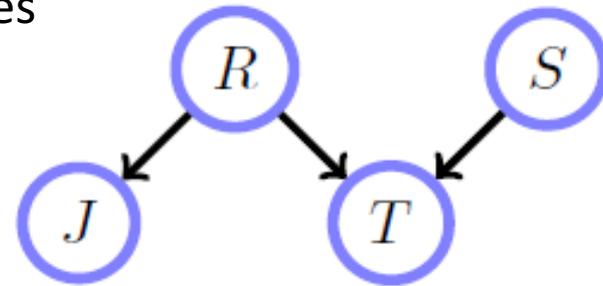
Directed Graphical Models

- What kinds of information can and cannot be encoded in the graph?
- DGM
 - Can encode only simplifying assumptions about which variables are conditionally independent from each other
 - Only defines which variables they are allowed to take in as arguments.
- But, DGM
 - Cannot encode the assumption that Bob's personal running time is independent from all other factors
 - Where The conditional distribution is now a slightly more complicated formula using only $k - 1$ parameters
 - Does not place any constraint on how we define our conditional distributions

Wet grass example

One morning Tracey leaves her house and realises that her grass is wet. Is it due to overnight rain or did she forget to turn o the sprinkler last night? Next she notices that the grass of her neighbour, Jack, is also wet.

This **explains away** to some extent the possibility that her sprinkler was left on, and she concludes therefore that it has probably been raining



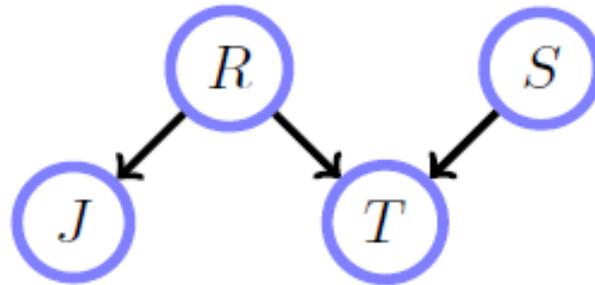
$R \in \{0, 1\}$ $R = 1$ means that it has been raining, and 0 otherwise

$S \in \{0, 1\}$ $S = 1$ means that Tracey has forgotten to turn off the sprinkler, and 0 otherwise

$J \in \{0, 1\}$ $J = 1$ means that Jack's grass is wet, and 0 otherwise

$T \in \{0, 1\}$ $T = 1$ means that Tracey's Grass is wet, and 0 otherwise

Wet grass example: Conditional independence



Conditional independence

$$\begin{array}{ll} P(T|J,R,S)=P(T|R,S) & P(J|R,S)=P(J|R) \\ P(J|R,S)=P(J|R) & P(R|S)=P(R) \end{array}$$



$$P(T,J,R,S)=P(T|R,S)P(J|R)P(R)P(S)$$

Wet grass example: Conditional Probability Table & Inference

$$p(R = 1) = 0.2, p(S = 1) = 0.1.$$

$$p(J = 1 | R = 1) = 1, p(J = 1 | R = 0) = 0.2$$

$$p(T = 1 | R = 1, S = 0) = 1, p(T = 1 | R = 1, S = 1) = 1,$$

$$p(T = 1 | R = 0, S = 1) = 0.9$$

$$p(T = 1 | R = 0, S = 0) = 0.$$

↓ Inference

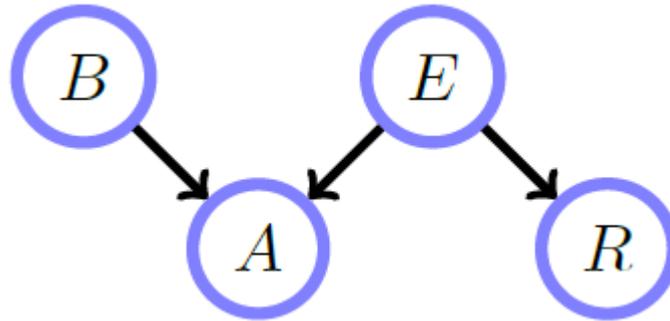
$$\begin{aligned} p(S = 1 | T = 1) &= \frac{p(S = 1, T = 1)}{p(T = 1)} = \frac{\sum_{J,R} p(T = 1, J, R, S = 1)}{\sum_{J,R,S} p(T = 1, J, R, S)} \\ &= \frac{\sum_{J,R} p(J|R)p(T = 1|R, S = 1)p(R)p(S = 1)}{\sum_{J,R,S} p(J|R)p(T = 1|R, S)p(R)p(S)} \\ &= \frac{\sum_R p(T = 1|R, S = 1)p(R)p(S = 1)}{\sum_{R,S} p(T = 1|R, S)p(R)p(S)} \\ &= \frac{0.9 \times 0.8 \times 0.1 + 1 \times 0.2 \times 0.1}{0.9 \times 0.8 \times 0.1 + 1 \times 0.2 \times 0.1 + 0 \times 0.8 \times 0.9 + 1 \times 0.2 \times 0.9} = 0.3382 \end{aligned}$$

Wet grass example: Inference

$$\begin{aligned} p(S = 1|T = 1, J = 1) &= \frac{p(S = 1, T = 1, J = 1)}{p(T = 1, J = 1)} \\ &= \frac{\sum_R p(T = 1, J = 1, R, S = 1)}{\sum_{R,S} p(T = 1, J = 1, R, S)} \\ &= \frac{\sum_R p(J = 1|R)p(T = 1|R, S = 1)p(R)p(S = 1)}{\sum_{R,S} p(J = 1|R)p(T = 1|R, S)p(R)p(S)} \\ &= \frac{0.0344}{0.2144} = 0.1604 \end{aligned}$$

The probability that the sprinkler is on, given the extra evidence that Jack's grass is wet, is lower than the probability that the grass is wet given only that Tracey's grass is wet. This occurs since the fact that Jack's grass is also wet increases the chance that the rain has played a role in making Tracey's grass wet.

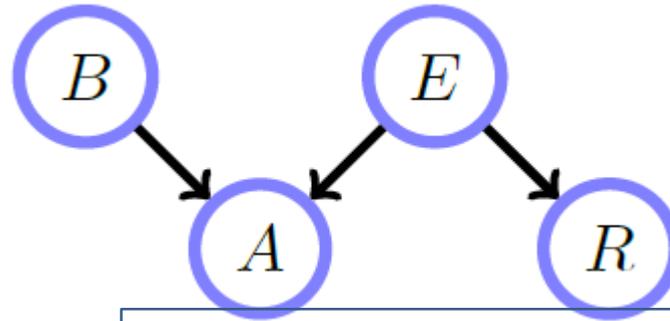
Burglar model example



Sally comes home to find that the burglar alarm is sounding ($A = 1$). Has she been burgled ($B = 1$), or was the alarm triggered by an earthquake ($E = 1$)?

She turns the car radio on for news of earthquakes and finds that the radio broadcasts an earthquake alert ($R = 1$).

Burglar model example



Radio = 1	Earthquake
1	1
0	0

$$p(B, E, A, R) = p(A|B, E, R)p(B, E, R)$$

$$p(B, E, A, R) = p(A|B, E, R)p(R|B, E)p(E|B)p(B)$$

$$p(B, E, A, R) = p(A|B, E)p(R|E)p(E)p(B)$$

Alarm = 1	Burglar	Earthquake
0.9999	1	1
0.99	1	0
0.99	0	1
0.0001	0	0

$$p(B = 1) = 0.01 \text{ and } p(E = 1) = 0.000001.$$

Burglar model example

- Initial Evidence: The Alarm is sounding

$$\begin{aligned} p(B = 1|A = 1) &= \frac{\sum_{E,R} p(B = 1, E, A = 1, R)}{\sum_{B,E,R} p(B, E, A = 1, R)} \\ &= \frac{\sum_{E,R} p(A = 1|B = 1, E)p(B = 1)p(E)p(R|E)}{\sum_{B,E,R} p(A = 1|B, E)p(B)p(E)p(R|E)} \approx 0.99 \end{aligned}$$

- Additional Evidence: The Radio broadcasts an Earthquake warning

$$p(B = 1|A = 1, R = 1) \approx 0.01$$



the Earthquake 'explains away' to an extent the fact that the Alarm is ringing.

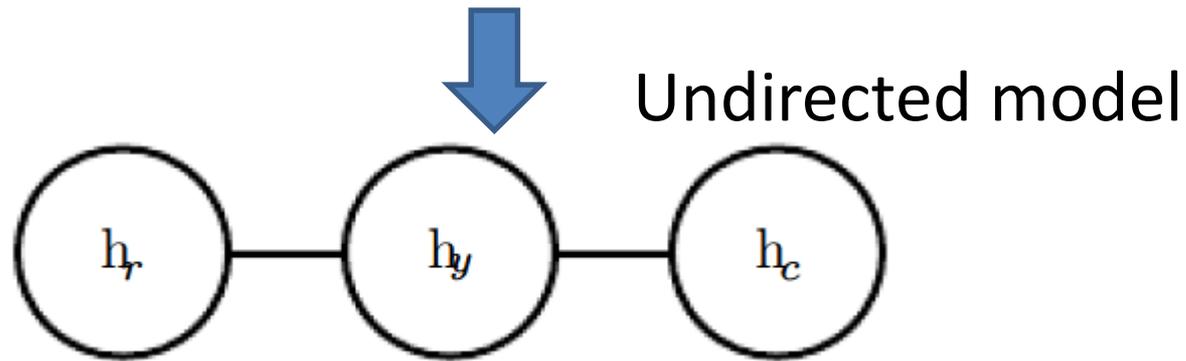
Undirected Graphical Models

- known as *Markov random fields* (MRFs) or *Markov networks*
 - Not all situations we might want to model have such a clear direction to their interactions
 - When the interactions seem to have no intrinsic direction, or to operate in both directions, it may be more appropriate to use an undirected model

Cold Spreading Example

Suppose that we want to model a distribution over three binary variables: whether or not you are sick, whether or not your coworker is sick, and whether or not your roommate is sick

No clean, uni-directional narrative on which to base the model.



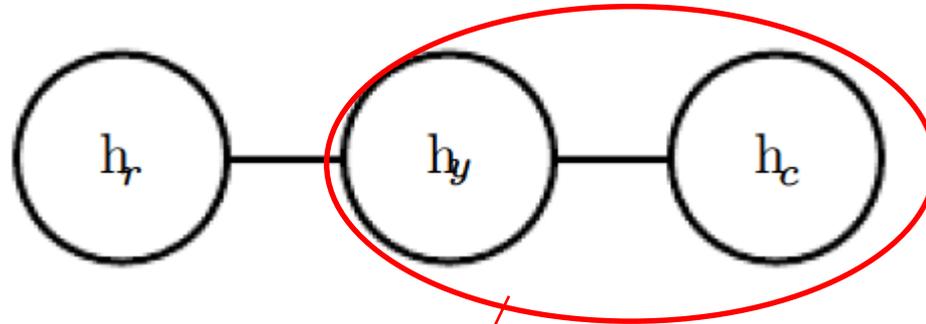
An undirected graph representing how your roommate's health h_r , your health h_y , and your work colleague's health h_c affect each other.

Undirected Graphical Models

- a structured probabilistic model defined on an undirected graph G .
- For each clique C in the graph, a *factor* $\phi(C)$ (also called a *clique potential*) measures the affinity of the variables in that clique for being in each of their possible joint states.
- The factors define *unnormalized probability distribution*

$$\tilde{p}(\mathbf{x}) = \prod_{C \in \mathcal{G}} \phi(C)$$

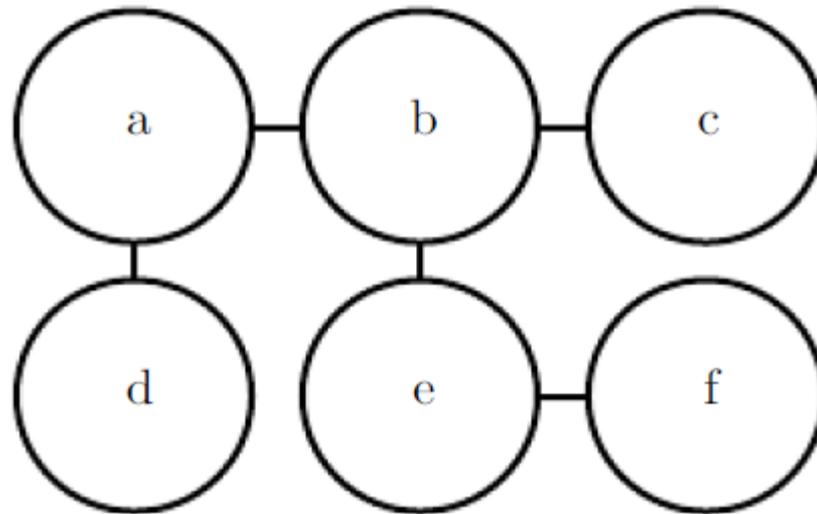
Cold Spreading Example



	$h_y = 0$	$h_y = 1$
$h_c = 0$	2	1
$h_c = 1$	1	10

a state of 0 indicates poor health

Undirected Graphical Models



$$p(a, b, c, d, e, f)$$



$$\frac{1}{Z} \phi_{a,b}(a, b) \phi_{b,c}(b, c) \phi_{a,d}(a, d) \phi_{b,e}(b, e) \phi_{e,f}(e, f)$$

The Partition Function

- The normalized probability distribution

$$p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x}) \quad Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$$

the partition function

- Computing Z exactly is usually intractable → we must resort to approximations
- There are often cases that Z does not exist
 - When the integral of \tilde{p} over their domain diverges.

- $\phi(x) = x^2$

$$Z = \int x^2 dx$$

Directed modeling vs. Undirected modeling

- Directed models: defined directly in terms of probability distributions from the start
- Undirected models: **defined more loosely** by ϕ functions that are then converted into probability distributions.
 - **the domain of each of the variables** has dramatic effect on the kind of probability distribution

- E.g.) $\phi^{(i)}(x_i) = \exp(b_i x_i)$.

$$\mathbf{x} \in \mathbb{R}^n$$



No prob
distribution

$$\mathbf{x} \in \{0, 1\}^n$$



$$p(x_i = 1) = \text{sigmoid}(b_i)$$

$$\mathbf{x} \in$$

$$\{[1, 0, \dots, 0], [0, 1, \dots, 0], \dots, [0, 0, \dots, 1]\}$$



$$p(\mathbf{x}) = \text{softmax}(\mathbf{b})$$

Energy-Based Models

- Undirected models assume $\forall \mathbf{x}, \tilde{p}(\mathbf{x}) > 0$

Energy-based model (EBM)

$$\tilde{p}(\mathbf{x}) = \exp(-E(\mathbf{x}))$$

energy function



- EBM is an example of a *Boltzmann distribution*
 - many energy-based models are called *Boltzmann machines*

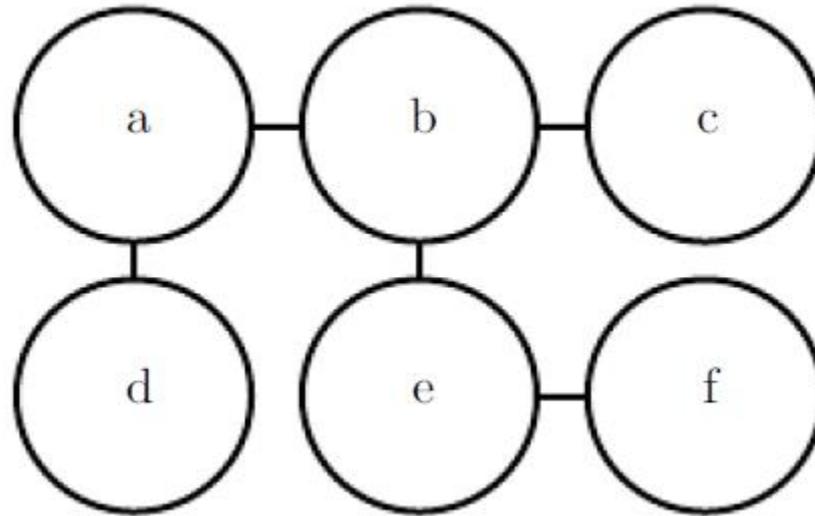
Boltzmann Machine

- A Boltzmann machine is a MN on binary variables $dom(x_i) = \{0, 1\}$ of the form

$$p(\mathbf{x}) = \frac{1}{Z(\mathbf{w}, b)} e^{\sum_{i < j} w_{ij} x_i x_j + \sum_i b_i x_i}$$

- the term Boltzmann machine is today most often used to designate models with latent variables
- Boltzmann machines without latent variables are more often called [Markov random fields or log-linear models](#).

Energy-Based Models



$$E(a, b, c, d, e, f)$$

$$E_{a,b}(a, b) + E_{b,c}(b, c) + E_{a,d}(a, d) + E_{b,e}(b, e) + E_{e,f}(e, f)$$

$$\phi_{a,b}(a, b) = \exp(-E(a, b))$$

Energy-Based Models with Latent Variables

- Many algorithms that operate on probabilistic models do not need to compute $p_{model}(\mathbf{x})$ but only $\log \tilde{p}_{model}(\mathbf{x})$.
- *Free energy*

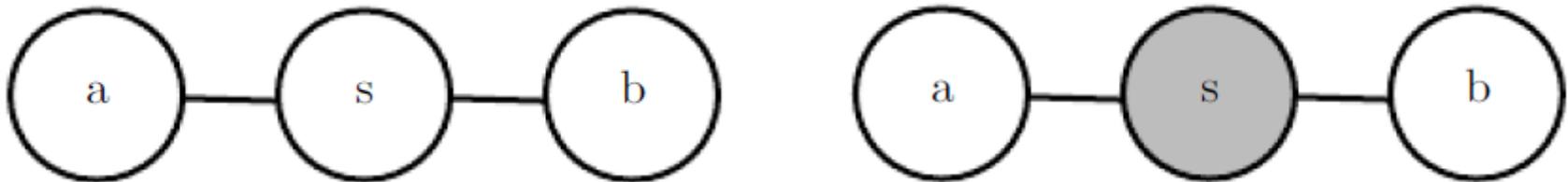
$$\mathcal{F}(\mathbf{x}) = -\log \sum_h \exp(-E(\mathbf{x}, \mathbf{h}))$$

Separation

which subsets of variables are conditionally independent from each other, given the values of other subsets of variables?

- a set of variables A is *separated* from another set of variables B given a third set of variables S if the graph structure implies that A is independent from B given S .

If two variables a and b are connected by a path involving only unobserved variables, then those variables are not separated

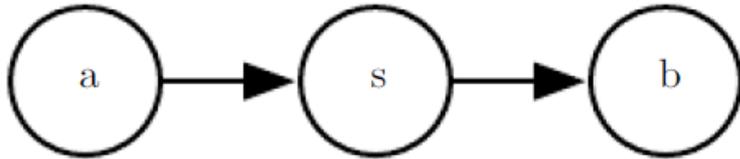


- paths involving only unobserved variables: **active**
- paths including an observed variable: **inactive**

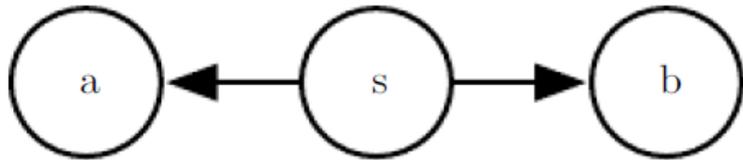
D-Separation

- The “d” stands for “dependence.”
- D-separation for directed graphs
- a set of variables A is d-separated from another set of variables B given a third set of variables S if the graph structure implies that A is independent from B given S .

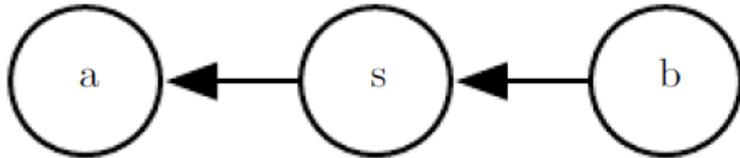
D-Separation



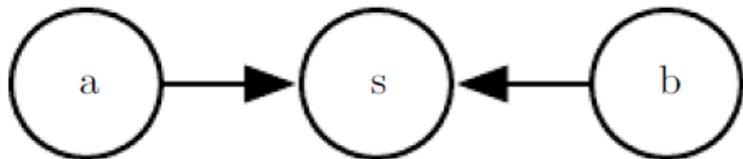
(a)



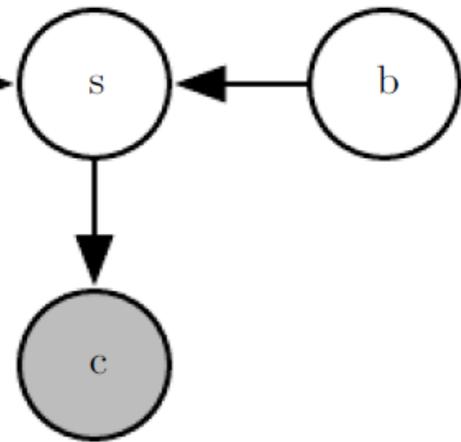
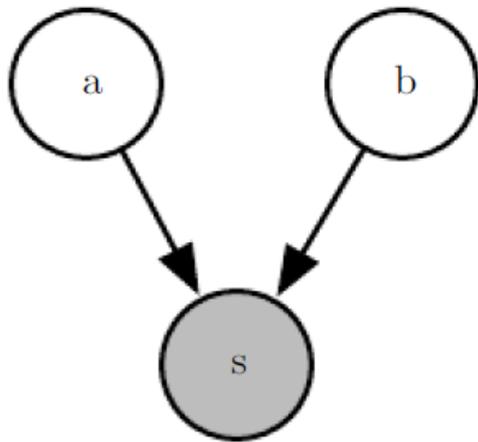
(b)



(c)

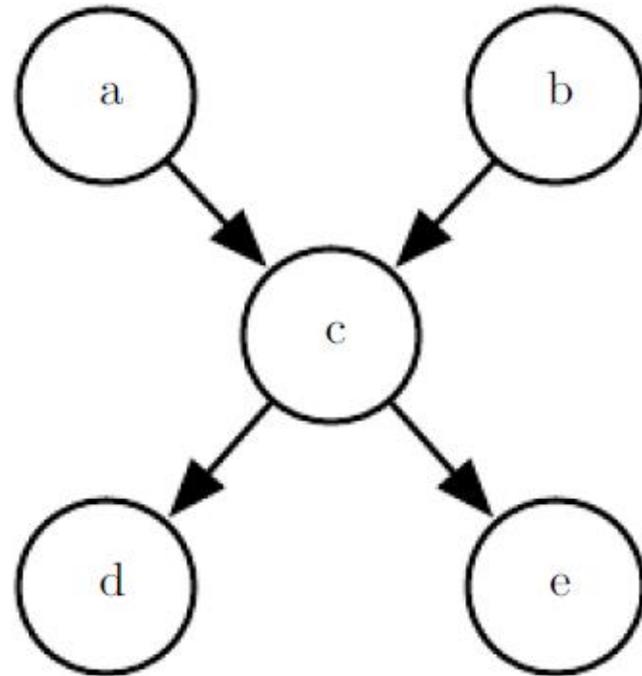


(d)



D-Separation

- a and b are d-separated given the empty set.
- a and e are d-separated given c.
- d and e are d-separated given c.



- a and b are not d-separated given c.
- a and b are not d-separated given d.

Separation and D-Separation

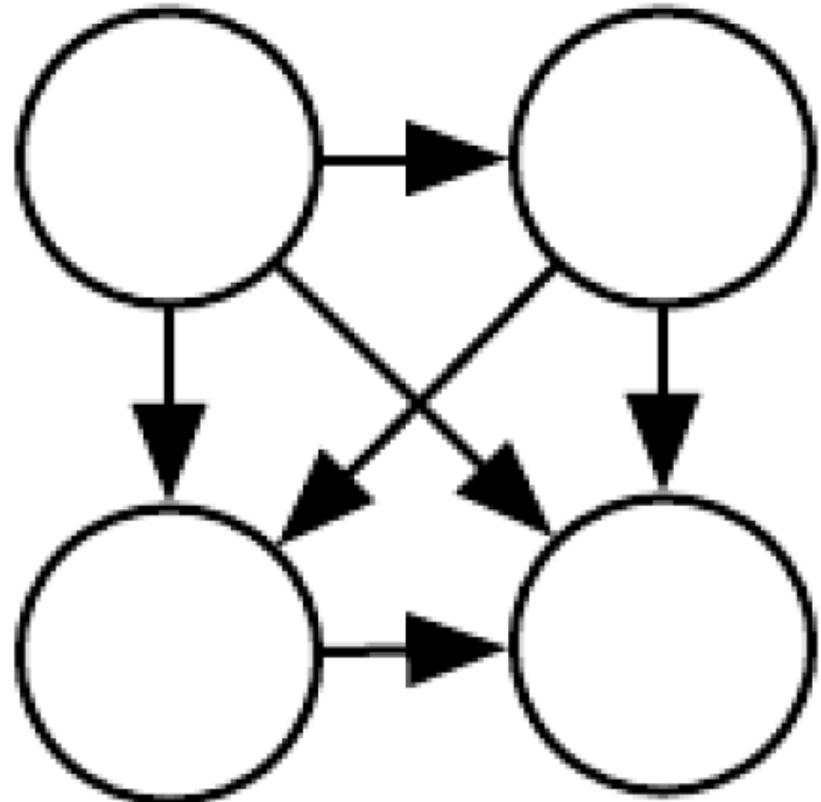
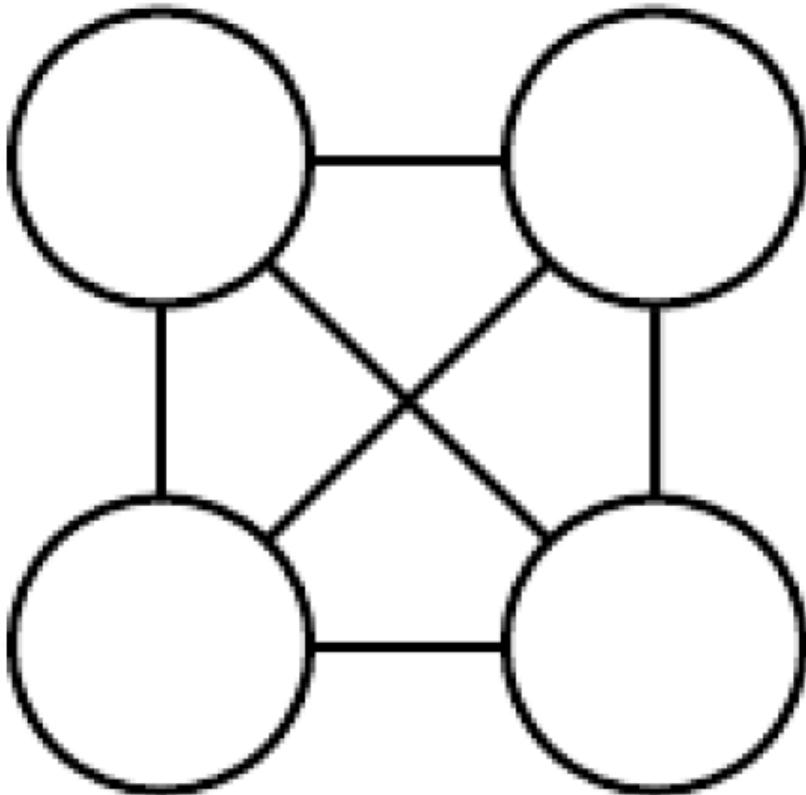
- Separation and d-separation tell us only about those **conditional independences that are implied by the graph**
- No requirement that the graph imply all independences that are present.
- ***Context-specific independences***: independences that are present dependent on the value of some variables in the network.
 - These independences are not possible to represent with existing graphical notation:

Separation and D-Separation

- In general, a graph will never imply that an independency exists when it does not
- However, a graph may fail to encode an independence

Converting between Undirected and Directed Graphs

Some models are most easily described using a directed graph, or most easily described using an undirected graph.



Converting Directed Models to Undirected Models

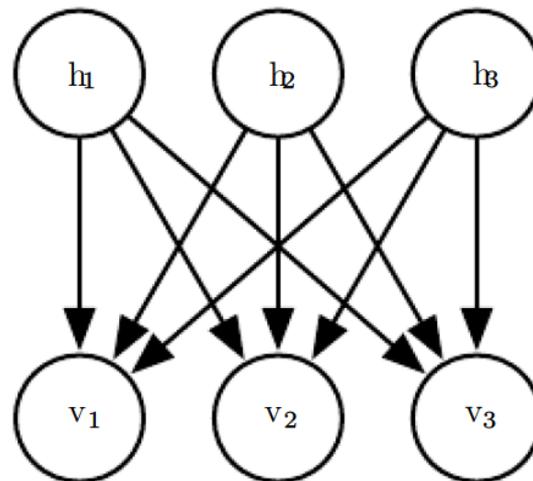
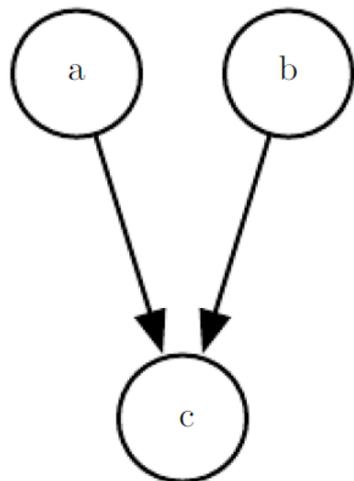
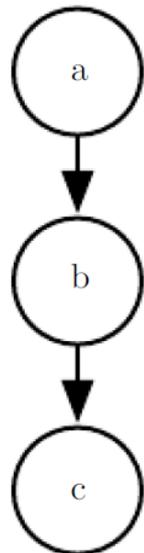
- *Immorality*: Directed models are able to use one specific kind of substructure that undirected models cannot represent perfectly
 - Occurs when **two random variables a and b are both parents of a third random variable c**, and there is **no edge directly connecting a and b** in either direction.



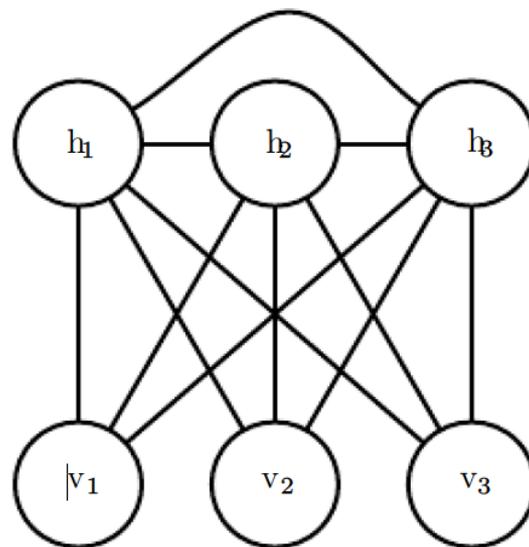
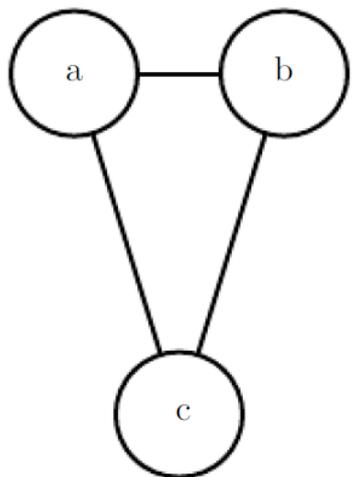
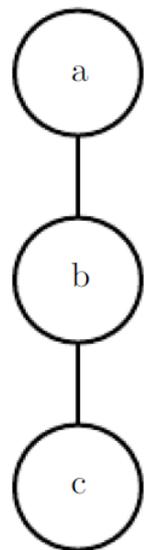
Moralization

Add an undirected edge connecting parents

Converting Directed Models to Undirected Models



Moralization



Moralized graphs

Converting undirected model to directed model

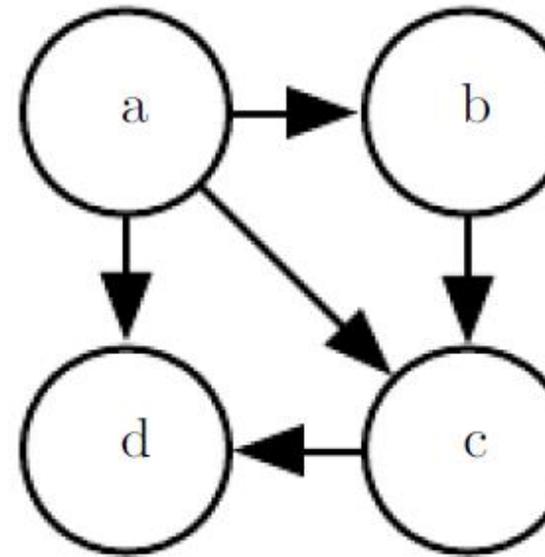
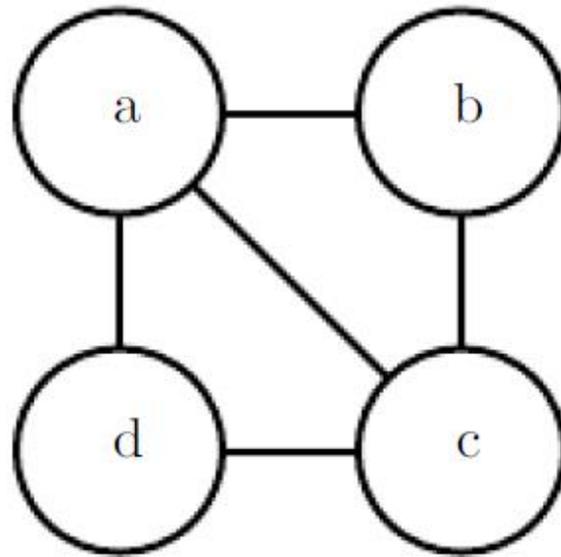
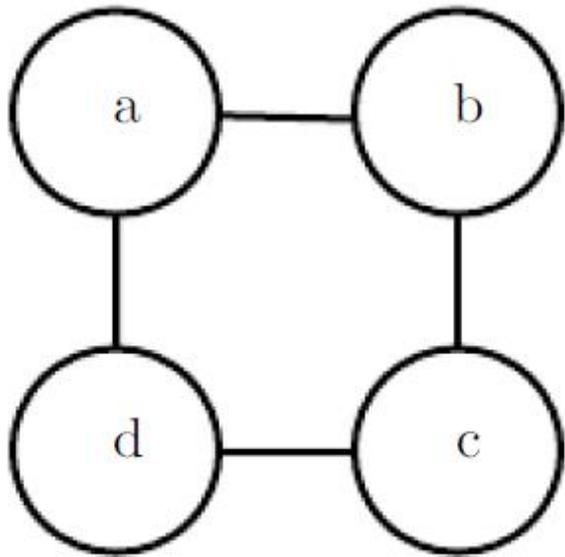
- a directed graph D cannot capture all of the conditional independences implied by an undirected graph U if U contains a *loop* of length greater than three, unless that loop also contains a *chord*.



Triangulation

- adding chords to U is known as a *chordal* or *triangulated graph*,

Converting undirected model to directed model

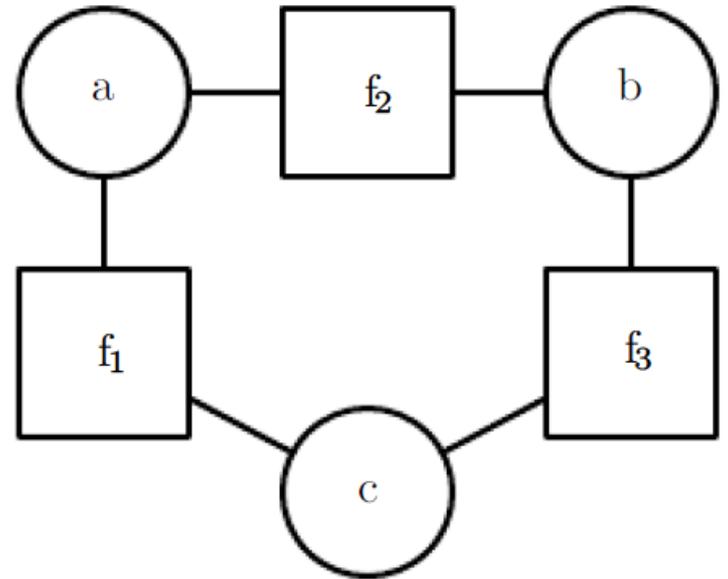
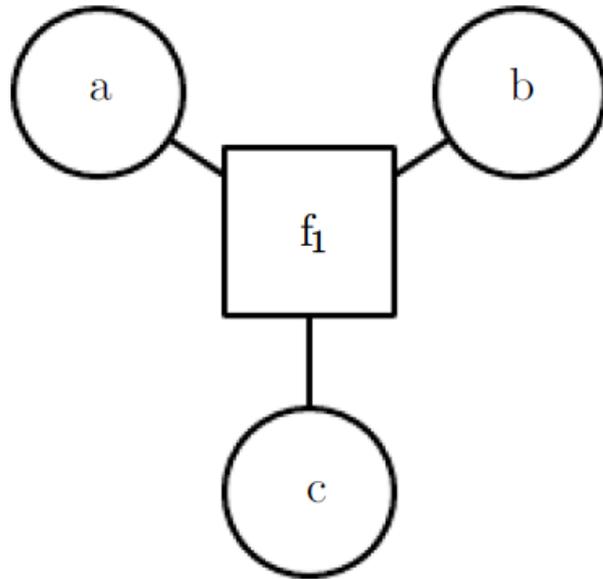
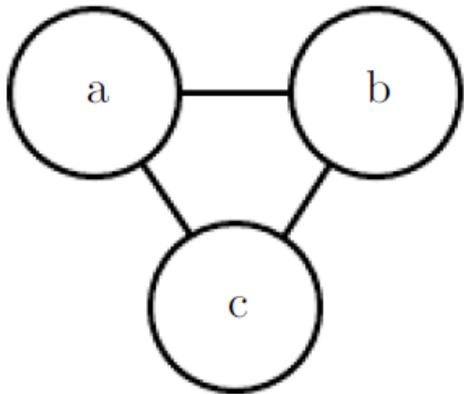


triangulated graph

Factor Graphs

- graphical representation of an undirected model that consists of a bipartite undirected graph
- **Resolve an ambiguity in** the graphical representation of standard undirected model syntax
- Variable nodes: drawn as circles
- Factor nodes: drawn as squares
 - correspond to the factors ϕ of the unnormalized probability

Factor Graphs



Sampling from Graphical Models

- *Ancestral sampling*
 - Sort the variables x_i in the graph into a topological ordering
 - Sample $x_1 \sim P(x_1)$,
 - sample $P(x_2 | Pa_G(x_2))$,
 - ...
 - Only applicable for directed graphical models

Sampling from Graphical Models

- *Gibbs sampling*
 - iteratively visit each variable x_i and draw a sample conditioned on all of the other variables, $P(x_i | x_{-i})$
 - Repeat the process and resample all n variables using the updated values of their neighbors
 - Asymptotically, after many repetitions, this process converges to sampling from the correct distribution
 - Can draw samples from an undirected graphical model

Advantages of Structured Modeling

- Dramatically reduce the cost of representing probability distributions as well as learning and inference
- Explicitly separate representation of knowledge from learning of knowledge or inference given existing knowledge

Learning about Dependencies

- A good generative model needs to accurately capture the distribution over the observed or **visible** variables v .
 - Can use **structured learning based** on greedy search
- In the context of deep learning, the approach most commonly used to model these dependencies is to introduce several latent or **hidden** variables, h .
 - Accomplish feature learning by learning latent variables

Inference and Approximate Inference

- In a latent variable model, we might want to extract features $E[\mathbf{h} \mid \mathbf{v}]$ describing the observed variables \mathbf{v} .
- Training based on ML

$$\log p(\mathbf{v}) = \mathbb{E}_{\mathbf{h} \sim p(\mathbf{h}|\mathbf{v})} [\log p(\mathbf{h}, \mathbf{v}) - \log p(\mathbf{h} \mid \mathbf{v})]$$

Inference problems

we must predict the value of some variables given other variables, or predict the probability distribution over some variables given the value of other variables



Intractable \rightarrow *Approximate inference*

Deep Learning Approach to Structured Probabilistic Models

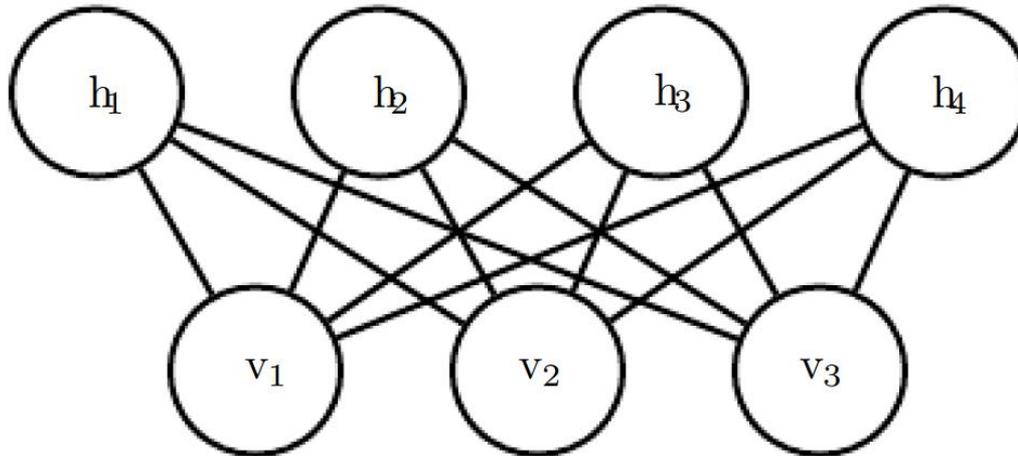
- Does not always involve especially deep graphical models
- Essentially always makes use of the idea of distributed representations.
 - Typically have more latent variables than observed variables.
 - By contrast, traditional graphical models usually contain mostly variables that are at least occasionally observed

Deep Learning Approach to Structured Probabilistic Models

- Does not intend for the latent variables to take on any specific semantics ahead of time
 - But, in the context of traditional graphical models, they are often designed with some specific semantics in mind
 - E.g. the topic of a document,
- Typically have large groups of units that are all connected to other groups of units
 - The interactions between two groups may be described by a single matrix.
 - But, traditional graphical models have very few connections and the choice of connections for each variable may be individually designed.

Example: RBM

- The typical deep learning approach to graphical models
 - Representation learning accomplished via **layers of latent variables**, combined with **efficient interactions** between layers parametrized by matrices



$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h}$$

$$p(\mathbf{h} \mid \mathbf{v}) = \prod_i p(h_i \mid \mathbf{v})$$

$$p(\mathbf{v} \mid \mathbf{h}) = \prod_i p(v_i \mid \mathbf{h})$$

Example: RBM

$$P(h_i = 1 \mid \mathbf{v}) = \sigma \left(\mathbf{v}^\top \mathbf{W}_{:,i} + b_i \right)$$

$$P(h_i = 0 \mid \mathbf{v}) = 1 - \sigma \left(\mathbf{v}^\top \mathbf{W}_{:,i} + b_i \right)$$

➡ Allow for efficient *block Gibbs* sampling

$$\frac{\partial}{\partial W_{i,j}} E(\mathbf{v}, \mathbf{h}) = -v_i h_j$$

Efficient Gibbs sampling and efficient derivatives make training convenient.

- Training the model induces a representation \mathbf{h} of the data \mathbf{v}

$\mathbb{E}_{\mathbf{h} \sim p(\mathbf{h}|\mathbf{v})} [\mathbf{h}] \rightarrow$ A set of features to describe \mathbf{v}