

Using Graphs to Describe Model Structure

Sargur N. Srihari
srihari@cedar.buffalo.edu

Topics in Structured PGMs for Deep Learning

0. Overview

1. Challenge of Unstructured Modeling

2. Using graphs to describe model structure

3. Sampling from graphical models

4. Advantages of structured modeling

5. Learning about Dependencies

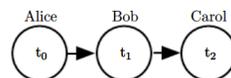
6. Inference and Approximate Inference

7. The deep learning approach to structured probabilistic models

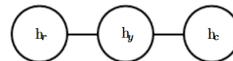
1. Ex: The Restricted Boltzmann machine

Topics in Graphs to Describe Model Structure

1. Directed Models



2. Undirected Models



3. The Partition Function

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

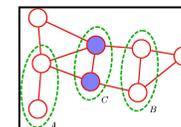
4. Energy-based Models

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

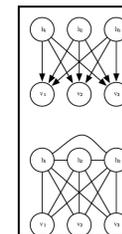
5. Restricted Boltzmann Machine

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

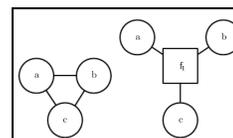
6. Separation and D-separation



7. Converting between Undirected and Directed Graphs



8. Factor Graphs



Graphs to describe model structure

- Model structure is described using graphs:
 - Each node represents a random variable
 - Each edge represents a direct interaction
 - These direct interactions imply other indirect interactions
 - But only direct interactions need be explicitly modeled

Types of graphical models

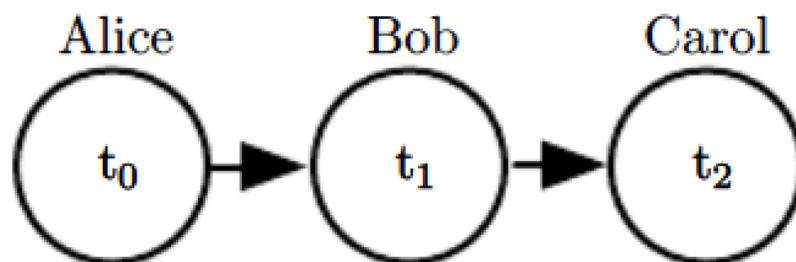
- More than one way to describe interactions in a probability distribution using a graph
- Graphical models divided into two categories
 - Models based on directed acyclic graphs
 - Models based on undirected graphs

1. Directed Models

- One type of structured probabilistic model is the directed graphical model
- Also known as
 - Belief network
 - Bayesian network (BN)
- The term Bayesian is used since the probabilities can be judgmental
 - Represent degrees of belief rather than frequencies of events

Example of Directed Graphical Model

- Relay race example



- Bob's finishing time t_1
 - depends on Alice's finishing time t_0
- Carol's finishing time t_2
 - depends on Bob's finishing time t_1

Meaning of directed edges

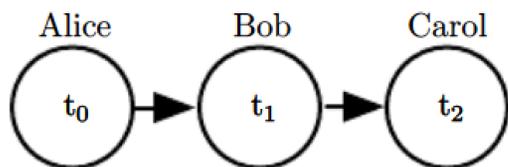
- Arrow from a to b means
 - we define a CPD over b via a conditional distribution with a as one of the variables on the right side of the conditional bar
- Distribution over b depends on the value of a

Formal directed graphical model

- Defined by a directed graphical acyclic graph G
 - Vertices are random variables in \mathbf{x}
 - A set of local CPDs $p(x_i | Pa_G(x_i))$
 - where $Pa_G(x_i)$ gives the parents of x_i in G
 - Distribution given by

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

- In the relay race example



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

Savings achieved by directed model

- If t_0, t_1, t_2 take 100 discrete values each
 - Full table has 999,999 values
 - CPDs of BN need only 18,999 values
- Cost of n discrete variables each with k values
 - Single table is $O(k^n)$
 - Directed model (BN):
 - If m is max. no. of variables in a CPD then CPDs need $O(k^m)$ values
- So long as each variable has few parents in graph, distribution represented by few parameters

2. Undirected models

- BNs give us a language to describe structured probabilistic models
- Another language is: undirected models
 - Also: Markov Random Fields, Markov Nets (MNs)
 - Use graphs whose edges are undirected
- BNs are useful when there is clear directionality
- When interactions have no clear directionality, more appropriate to use an undirected graph
 - Often when we know causality and causality flows in one direction

Models without clear direction

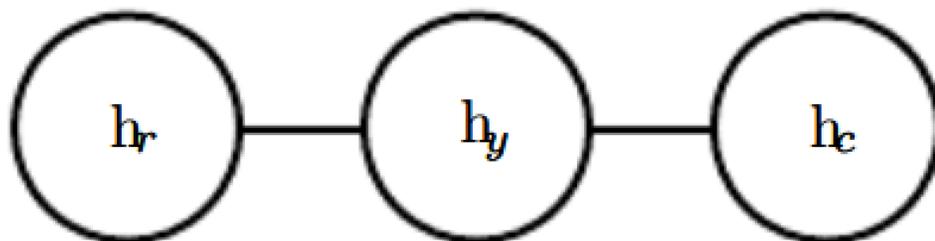
- When interactions have no clear directionality, or operate in both directions, it is appropriate to use an undirected model
- Example with three binary variables described next

Ex: Undirected Model for Health

- Three binary variables:
 - Whether you are sick or not, h_y
 - Whether your coworker is sick or not, h_c
 - Whether your roommate is sick or not, h_r
- Coworker and roommate do not know each other, so will not infect the other
 - Event is so rare we do not model it
- No clear directionality either
- This motivates using an undirected model

The health undirected graph

- You and your room-mate may infect each other
- You and your colleague may do the same
- Room-mate and colleague don't know each other they can only get infected through you



Undirected graph definition

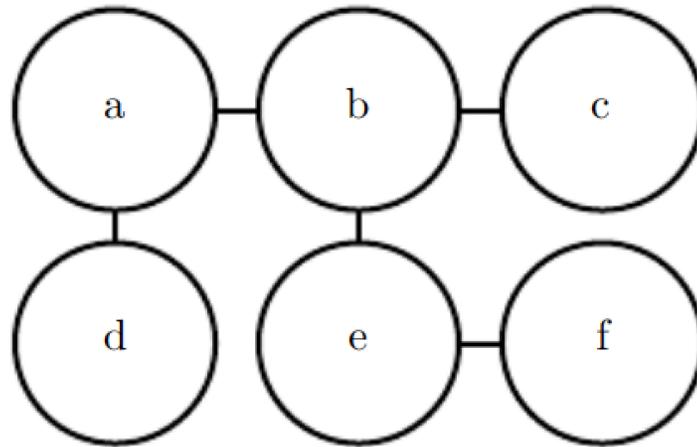
- If variables directly interact then nodes are connected
 - Edge has no arrow and has no CPD
- An undirected PGM is defined on a graph G
 - For clique C , a factor $\phi(C)$, or clique potential, measures affinity for being in each of their joint states
 - A clique is a subset of nodes all connected to each other
 - Together they define an unnormalized distribution

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

Efficiency of Unnormalized Distribution

- Unnormalized distribution is efficient to work with when cliques are small
- It encodes that states with higher affinity $\phi(C)$ are more likely
- Since there is little structure to the definition of the cliques, there is no guarantee that multiplying them together will yield a probability distribution

Reading factorization information from an undirected graph

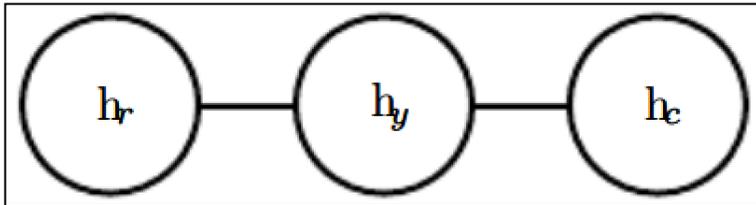


- Graph (with five cliques) implies

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

- for an appropriate choice of the ϕ functions
- Example of clique potentials shown next

Ex: Clique potential



h_y = health of you
 h_r = health of roommate
 h_c = health of colleague

- One clique is between h_y and h_c
 - Factor for this clique can be defined by a table

$\phi(h_y, h_c)$	$h_y = 0$	$h_y = 1$
$h_c = 0$	2	1
$h_c = 1$	1	10

A state of 1 indicates good health

While state of 0 indicates poor health

Both are usually healthy, so the corresponding state has highest affinity

State of only one being being sick has lowest affinity

State of both being sick has higher affinity than one being sick

- Similar factor needed for the other clique between h_y and h_c

3. The partition function

- Unnormalized distribution is
 - guaranteed to be non-negative everywhere
 - not guaranteed to sum to 1

- Normalized (Gibbs) distribution

$$p(\mathbf{x}) = \frac{1}{Z} \hat{p}(\mathbf{x}) \qquad \tilde{p}(\mathbf{x}) = \prod_{C \in G} \phi(C)$$

- Where Z causes the distribution to sum to one

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

- Z is a constant when the ϕ functions are constant
- If ϕ has parameters then Z is a function of those parameters, commonly written without arguments
- Known as the partition function in statistical physics

Intractability of Z

- Since Z is an integral or sum over all possible values of \mathbf{x} it is intractable to compute
- In order to compute a normalized probability of an undirected model:
 - Model structure and definitions of ϕ functions must be conducive to computing Z efficiently
 - In deep learning applications Z is intractable and we must resort to approximations

Choice of factors

- For some factors, Z does not exist!

1. If scalar $x \in \mathbb{R}$ with a single clique potential $\phi(x) = x^2$

$$\boxed{Z = \int x^2 dx} \quad \textcircled{x}$$

- This integral diverges

– Hence there is no probability distribution for this choice

2. Parameters of ϕ determine if distribution exists

- For $\phi(x; \beta) = \exp(-\beta x^2)$, β determines whether Z exists

– Positive β defines a Gaussian over x

– Other values of β make ϕ impossible to normalize

Key difference between BN & MN

- Directed models are:
 - Defined directly in terms of probability distributions
 - From the start
- Undirected models are:
 - Defined loosely in terms of ϕ functions that are then converted into probability distributions
 - This changes intuitions to work with these models
 - Key idea in working with MNs:
 - Domain of variables has a dramatic effect on kind of probability distributions a given set of ϕ functions corresponds to
 - We will see how we can define distributions for different domains

What distribution does an MN give?

- Consider an n -dimensional random variable $\mathbf{x} = \{x_i\}_{i=1, \dots, n}$
- And an undirected model parameterized with biases \mathbf{b}
 - Suppose we have one clique for each x_i : $\phi^{(i)}(x_i) = \exp(b_i x_i)$



- What kind of probability distribution is modeled?

$$p(\mathbf{x}) = \frac{1}{Z} \hat{p}(\mathbf{x})$$

$$\tilde{p}(\mathbf{x}) = \prod_{C \in G} \phi^{(i)}(x_i) = \exp(b_1 x_1 + \dots + b_n x_n)$$

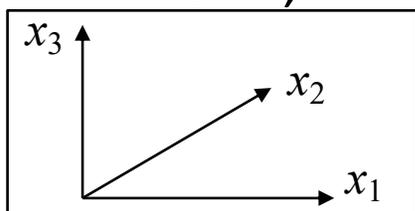
$$Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$$

- The answer is that we do not have enough information
- Because we have not specified the domain of \mathbf{x}

Domain of \mathbf{x} and resulting distributions

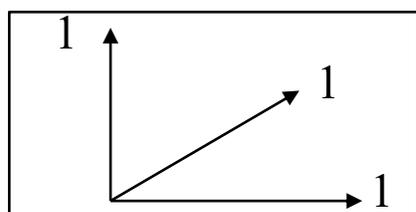
- We have n random variables, $\mathbf{x} = \{x_i\}_{i=1,\dots,n}$
- For each x_i : $\phi^{(i)}(x_i) = \exp(b_i x_i)$ $\tilde{p}(\mathbf{x}) = \prod_{C \in G} \phi^{(i)}(x_i) = \exp(b_1 x_1 + \dots + b_n x_n)$ $p(\mathbf{x}) = \frac{1}{Z} \hat{p}(\mathbf{x})$ $Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$
- Kind of probability distribution modeled depends on the domain

1. $\mathbf{x} \in \mathbb{R}^n$, an n -dimensional vector of real values



$Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$ diverges and no probability distribution exists

2. $\mathbf{x} \in \{0,1\}^n$, an n -dimensional vector of binary values



$$\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$p(\mathbf{x})$ factorizes into n independent distributions with $p(x_i=1) = \sigma(b_i)$, where each independent distribution is a binomial with parameter $\sigma(b_i)$

3. Domain of \mathbf{x} is set of basis vectors $\{[1,0,\dots,0],[0,1,\dots,0],\dots,[0,0,\dots,1]\}$

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$p(\mathbf{x}) = \text{softmax}(\mathbf{b})$
 So a large value of b_i reduces to $p(x_j=1)$ for $j \neq i$,
 i.e., multiclass

- By careful choice of domain of \mathbf{x} we can obtain complicated behavior from a simple set ϕ functions

4. Energy-based Models (EBMs)

- Many interesting theoretical results of undirected models depend on assumption that

$$\forall \mathbf{x}, \tilde{p}(\mathbf{x}) > 0 \quad \tilde{p}(\mathbf{x}) = \prod_{C \in \mathcal{G}} \phi(C)$$

- We can enforce this using an EBM where

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

- $E(\mathbf{x})$ is known as the Energy function
- Because $\exp(z) > 0 \forall z$, no energy function will result in a probability of zero for any \mathbf{x}
 - If we were to learn clique potentials directly we would need to impose constraints for minimum probability value
 - By learning the energy function we can use unconstrained optimization: probabilities can approach 0

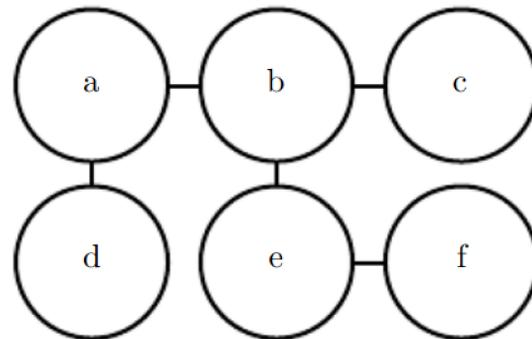
Boltzmann Machine Terminology

- Any distribution of the form $\tilde{p}(\mathbf{x}) = \exp(-E(\mathbf{x}))$
 - is referred to as a Boltzmann distribution
- For this reason, many energy-based models are referred to as Boltzmann machines
 - No consensus on when to call it a energy-based model or a Boltzmann Machine
 - Boltzmann machines first referred to only binary variables
 - Today mean-covariance restricted Boltzmann Machines deal with real-valued variables
 - Boltzmann Machines refer to models with latent variables and those without are referred to as MRFs or log-linear models

Cliques, factors and energy

- Cliques in the undirected graph correspond to factors in the unnormalized probability function
 - Cliques in undirected graph also correspond to different terms of an energy function
 - Because $\exp(a)\exp(b)=\exp(a+b)$ different cliques in undirected graph correspond to different terms of the energy function
 - i.e., energy-based model is a special Markov network
 - Exponentiation makes each term of the energy function correspond to a factor for a different clique
 - Reading the form of the energy function from an undirected graph is shown next

Graph and Corresponding Energy



- Graph, with five cliques, implies that

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

- We can obtain the ϕ functions in $p(a,b,c,d,e,f) = \frac{1}{Z} \phi_{a,b}(a,b) \phi_{b,c}(b,c) \phi_{a,d}(a,d) \phi_{b,e}(b,e) \phi_{e,f}(e,f)$

– By setting each ϕ to the exponential of the corresponding negative energy, e.g., if we set $\phi_{a,b}(a,b) = \exp(-E(a,b))$ we get

$$P(a,b,c,d,e,f) = (1/Z) \exp [-(E(a,b) + E(b,c) + E(a,d) + E(b,e) + E(e,f))]$$

- Negative sign serves no purpose

Role of negative sign in energy

- Negative sign in $\tilde{p}(\mathbf{x}) = \exp(-E(\mathbf{x}))$ serves no functional purpose from a ML perspective
- Sign could be incorporated into definition of Energy function
- It is for compatibility with physics literature
- Sometimes omit negative sign and refer to negative energy as *harmony*

Energy-based Model as Experts

- EBM with multiple terms is a *product of experts*
- Each term is a factor in the distribution
 - Term determines if a soft constraint is satisfied
 - An expert imposes one constraint
 - Concerns low-dimensional projection of random variables
- With multiplication of probabilities, experts together enforce a high-dimensional constraint

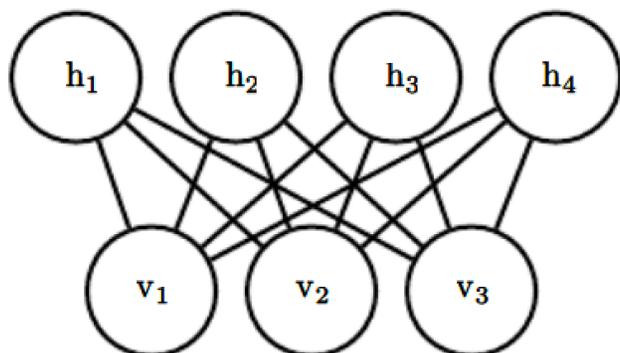
Free Energy instead of Probability

- Algorithms don't need $p_{\text{model}}(\mathbf{x})$ but only $\log \tilde{p}_{\text{model}}(\mathbf{x})$ where $\tilde{p}(\mathbf{x}) = \exp(-E(\mathbf{x}))$
- EBMs with hidden units \mathbf{h} use the negative of this quantity, called the *free energy*

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

– Deep learning prefers this formulation

RBM as a tractable Energy Model



Binary version of RBM

Observed layer: set of n_v binary r.v.s, \mathbf{v}

Latent or hidden set of n_h binary r.v.s \mathbf{h}

Its energy function is

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

where \mathbf{b} , \mathbf{c} and \mathbf{W} are learnable parameters

Interaction of units \mathbf{v} and \mathbf{h} is described by matrix \mathbf{W}

Joint-probability distribution of RBM is

$$P(\mathbf{v}, \mathbf{h}) = (1/Z) \exp(-E(\mathbf{v}, \mathbf{h}))$$

The energy function for an RBM is

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

Z is the partition function $Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))$

Since Z is intractable $P(\mathbf{v}, \mathbf{h})$ is also intractable

Tractability of RBM due to bipartite structure

Although $P(\mathbf{v}, \mathbf{h})$ is intractable, bipartite structure has special property that

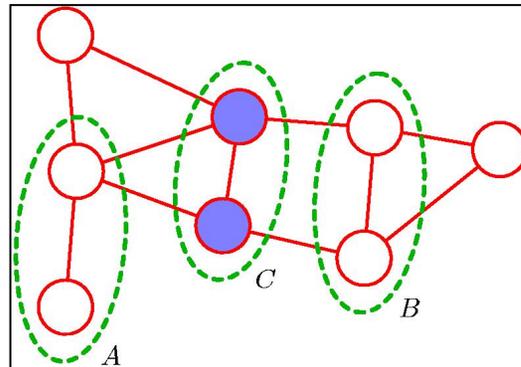
Conditionals $P(\mathbf{h}|\mathbf{v})$, $P(\mathbf{v}|\mathbf{h})$ are factorial & easily computed

5. Separation and D-Separation

- Edges in a directed graph tell which variables directly interact
- We often need to know which variables *indirectly* interact
 - Some of these can be enabled or disabled by observing other variables
 - More formally we need to know which variables are conditionally independent of each other given the values of other sets of variables

Separation in Undirected models

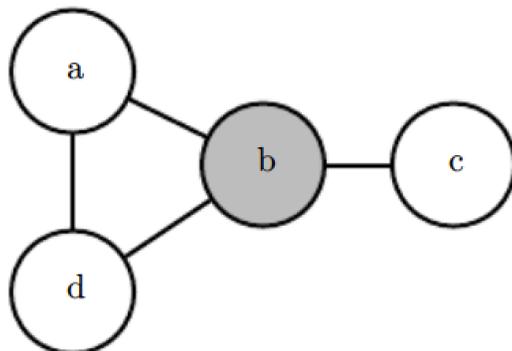
- Conditional independences are simple in MNs
 - Conditional independence is called **separation**
 - Variables A separated from variables B given variables C



- Graph implies A independent of B given C
- Path between variables a and b
 - All unobserved then not separated
 - No path exists then they are separated
 - All paths contain an observed variable then separated

Separation in undirected graphs

- **b** is shaded to indicate it is observed



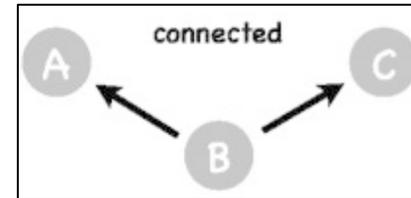
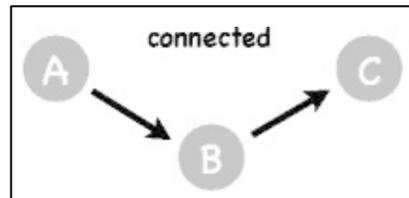
- **a-c are separated given b**
 - **b blocks path from a to c**
- **a-d not separated given b**
 - **There is an active path from a to d**

Separation in Directed Graphs

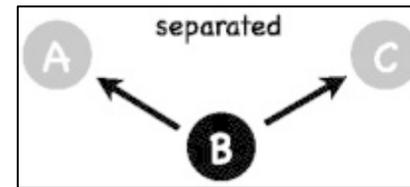
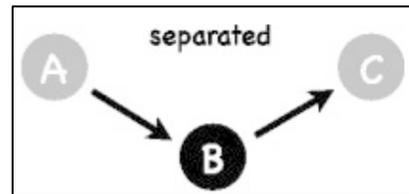
- In directed graphs, separation concepts are called D-separation
 - D stands for “dependence”
- D-separation is defined as separation for undirected graphs:
- A set of variables A is D-separated from a set of variables B given a third set of variables S if the graph structure implies that A is independent of B given S

Examining Active Paths

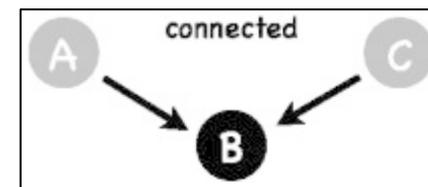
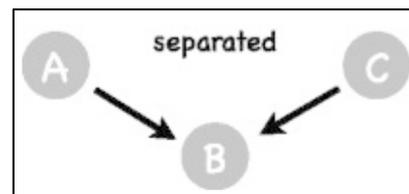
- Variables are dependent if active path exists



- D-separated if there is no path between them



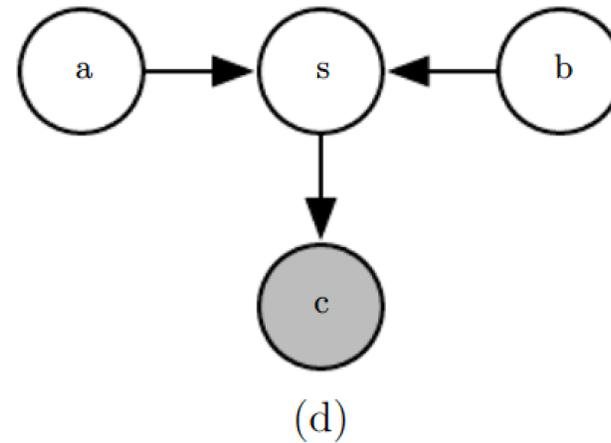
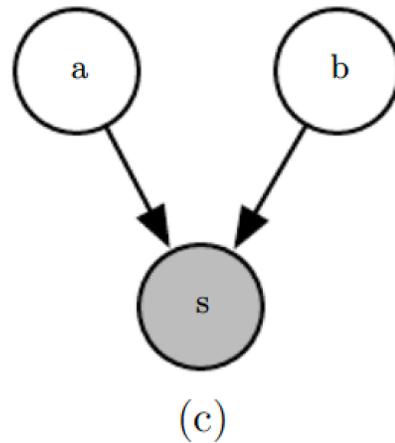
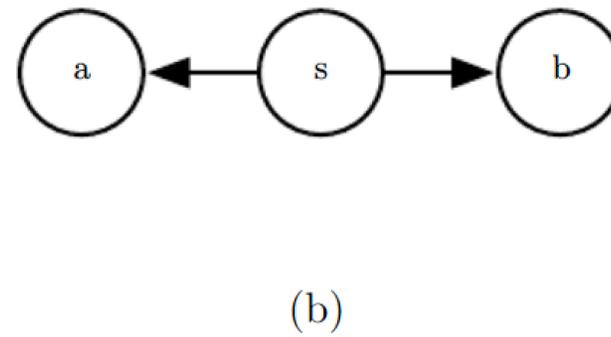
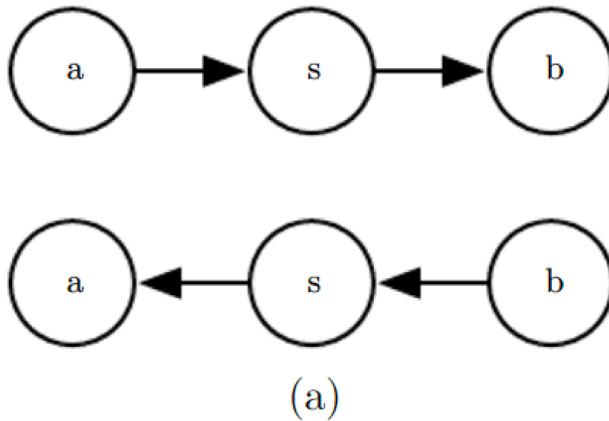
- Determining active path in BN is more complex



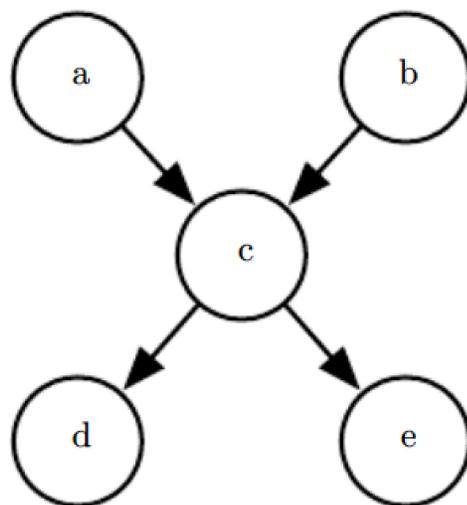
- Summary of active paths in a BN is given next

All active paths of length 2

Active paths between random variables a and b



Reading properties from a graph



- 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.

Undirected Graphs \leftrightarrow Directed Graphs

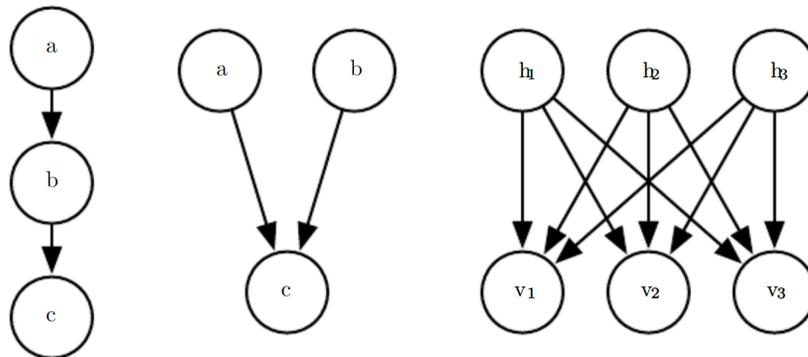
- ML models described as either undirected or directed
 - RBMs are undirected
 - Sparse coding are directed
- No model is inherently directed or undirected
 - Described easily using directed or undirected
 - Capture of independences or use fewer edges
- Strengths
 - Directed models provide straightforward sampling
 - Undirected models useful for approximate inference

Converting directed to undirected

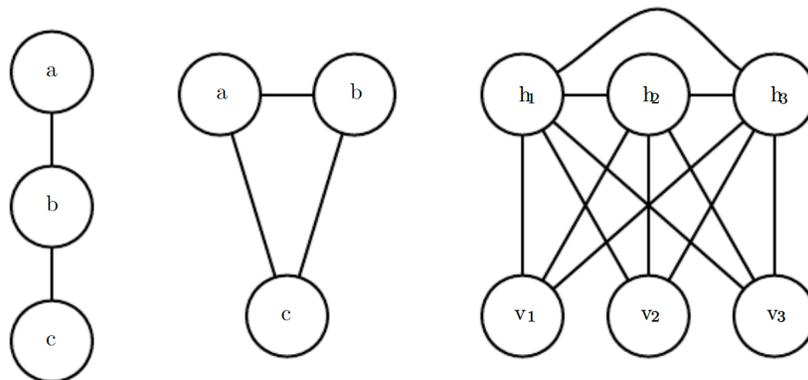
Converting directed models to undirected models via moralization

Resulting undirected model implies exactly the same set of independences and conditional independences

Directed Models

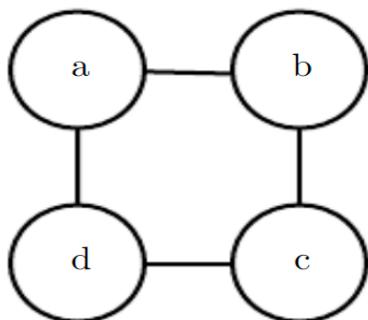


Undirected Models



Converting undirected to directed

This undirected model cannot be converted to a directed model because it has a loop of length four with no chords



Model encodes two different independences that no directed model can capture simultaneously:
 $a \perp c \mid \{b, d\}$ and $b \perp d \mid \{a, c\}$

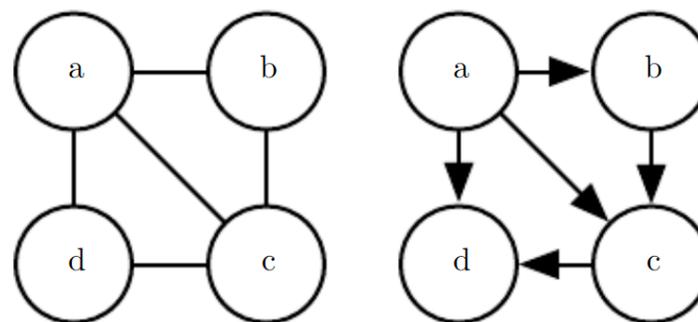
Triangulate the graph by ensuring that all loops of greater than length three have a chord.

Either add an edge connecting a and c or we can add an edge connecting b and d

Assign a direction to each edge.

When doing so, do not create any directed cycles.

Impose an ordering over nodes, and always point each edge from the node that comes earlier in the ordering to the node that comes later

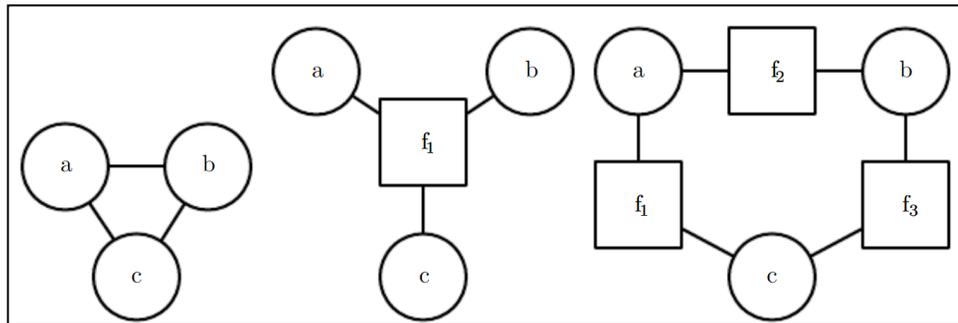


Factor Graphs

- Ambiguity in MNs

- Unclear if each clique actually has a corresponding factor whose scope encompasses the entire clique

- E.g., a clique containing three nodes may correspond to a factor over all three nodes, or may correspond to three factors that each contain only a pair of the nodes



- Representation, inference, and learning are cheaper in second factor graph than in the first, even though both require the same undirected graph to represent