

CSCI 699 - ProbGen

Probabilistic and Generative Models

Willie Neiswanger

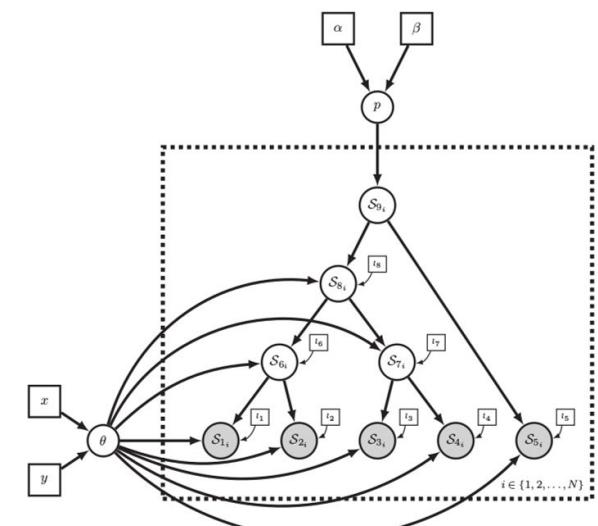
Lecture 2 - Introduction to PGMs

Plan for Today

Plan for Today

Lecture: PGMs – using ideas from both probability and graph theory to derive efficient machine learning algorithms

- Directed vs undirected graphical models: Bayesian networks and Markov random fields.
- Plate notation, generative process notation.
- Observed variables vs. latent variables vs. (hyper)parameters.
- Famous/classic Bayes Nets (HMMs, GMMs, LDA, VAE, etc.)
- Pros and Cons of Bayes Nets vs MRFs.
- Famous/classic MRFs (Ising Model, CRF).



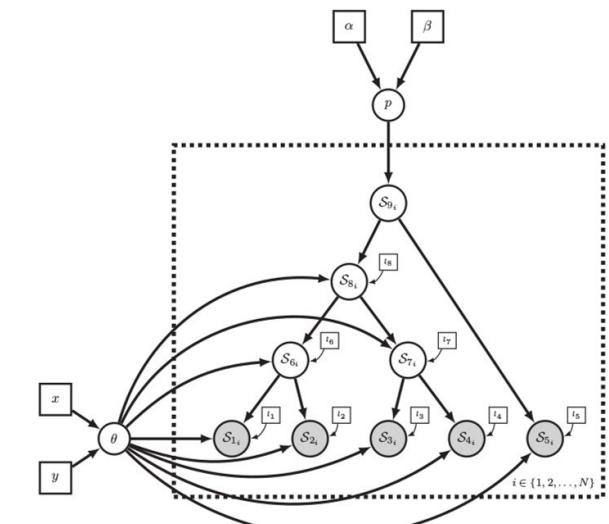
Plan for Today

Lecture: PGMs – using ideas from both probability and graph theory to derive efficient machine learning algorithms

- Directed vs undirected graphical models: Bayesian networks and Markov random fields.
- Plate notation, generative process notation.
- Observed variables vs. latent variables vs. (hyper)parameters.
- Famous/classic Bayes Nets (HMMs, GMMs, LDA, VAE, etc.)
- Pros and Cons of Bayes Nets vs MRFs.
- Famous/classic MRFs (Ising Model, CRF).

After:

- Class introductions and research interests, for group project.
- Compute resources.



Source: "Probabilistic Graphical Model Representation in Phylogenetics", Höhna 2014

Lecture 2 (January 24) – Intro to Probabilistic Graphical Models (PGMs)

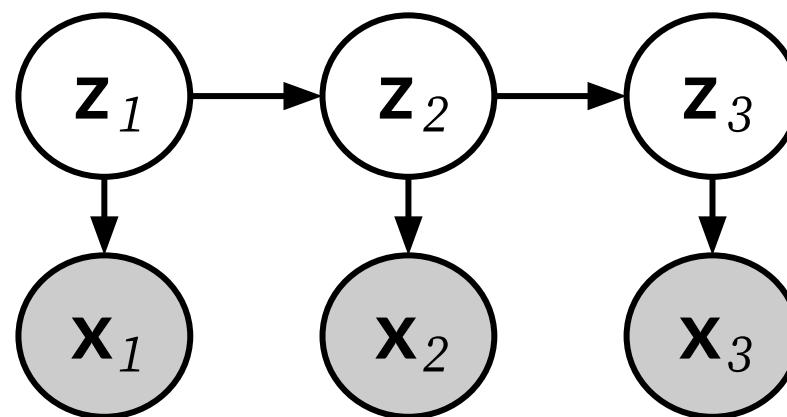
At a high level: two main “types” of PGMs.

Lecture 2 (January 24) – Intro to Probabilistic Graphical Models (PGMs)

At a high level: two main “types” of PGMs.

Bayesian Networks (directed)

E.g., Hidden Markov Model (HMM)

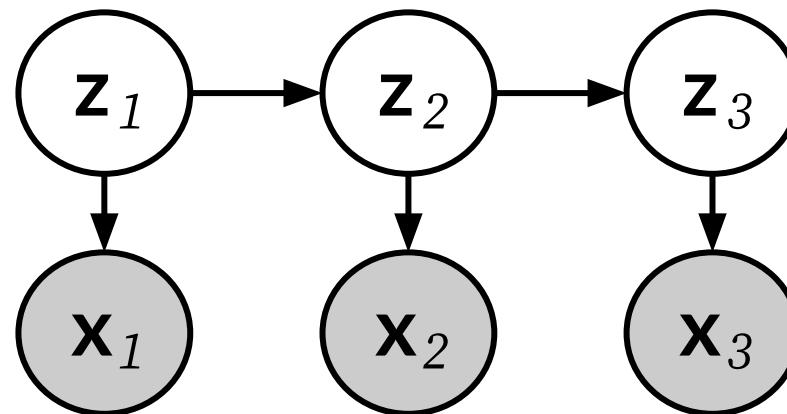


Lecture 2 (January 24) – Intro to Probabilistic Graphical Models (PGMs)

At a high level: two main “types” of PGMs.

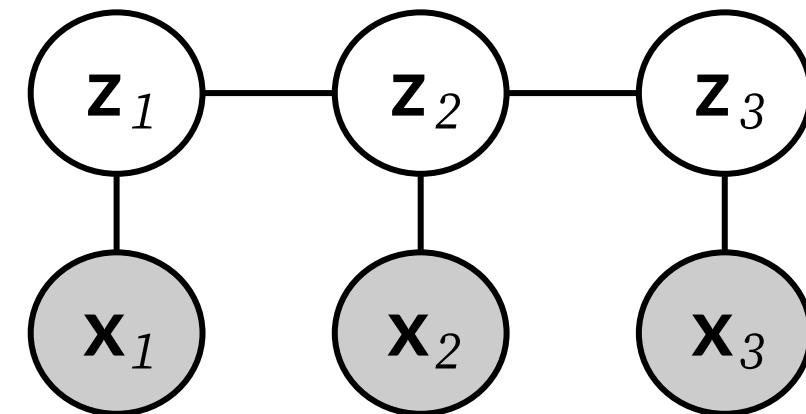
Bayesian Networks (directed)

E.g., Hidden Markov Model (HMM)



Markov Random Fields (undirected)

E.g., Conditional Random Field (CRF)

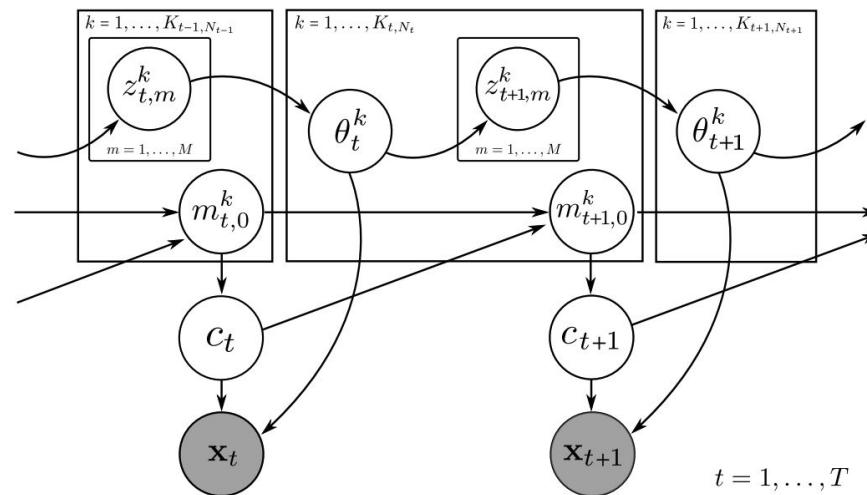


Lecture 2 (January 24) – Intro to Probabilistic Graphical Models (PGMs)

At a high level: two main “types” of PGMs.

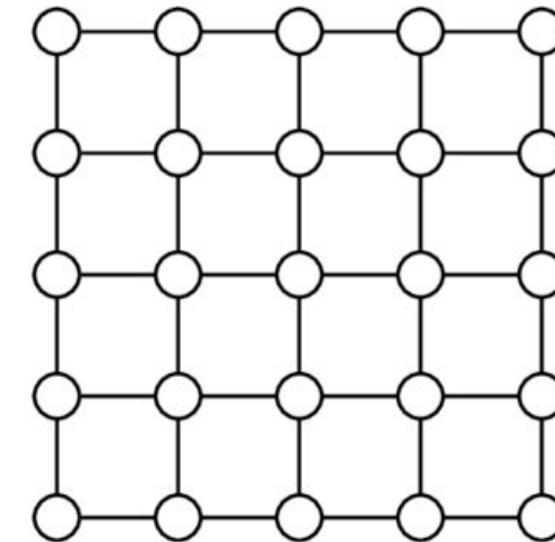
Bayesian Networks (directed)

E.g., Dependent Dirichlet Process Mixture of Experts (DDPMO)



Markov Random Fields (undirected)

E.g., Ising Model



Lecture 2 (January 24) – Intro to Probabilistic Graphical Models (PGMs)

Today, we will focus most of our time on **Bayesian networks**.

Lecture 2 (January 24) – Intro to Probabilistic Graphical Models (PGMs)

Today, we will focus most of our time on **Bayesian networks**.

They have certain advantages over Markov Random Fields, which we will discuss.

The converse is also true: MRFs have certain advantages.

Lecture 2 (January 24) – Intro to Probabilistic Graphical Models (PGMs)

Today, we will focus most of our time on **Bayesian networks**.

They have certain advantages over Markov Random Fields, which we will discuss.

The converse is also true: MRFs have certain advantages.

And also:

Bayesian networks are more relevant to some of the approximate inference algorithms and deep generative models we will cover a bit later in the class.

Directed PGMs: Bayesian Networks

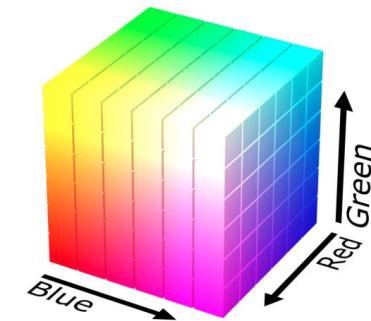
Image Example

Modeling a single pixel in an image

Image Example

Modeling a single pixel in an image \Rightarrow 3 discrete random variables:

- R - Red Channel. $\text{Val}(R) = \{0, \dots, 255\}$
- G - Green Channel. $\text{Val}(G) = \{0, \dots, 255\}$
- B - Blue Channel. $\text{Val}(B) = \{0, \dots, 255\}$

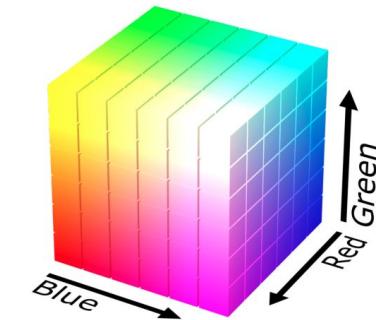


Source: Stefano Ermon, Deep Generative Models (CS236) Class

Image Example

Modeling a single pixel in an image \Rightarrow 3 discrete random variables:

- R - Red Channel. $\text{Val}(R) = \{0, \dots, 255\}$
- G - Green Channel. $\text{Val}(G) = \{0, \dots, 255\}$
- B - Blue Channel. $\text{Val}(B) = \{0, \dots, 255\}$



Source: Stefano Ermon, Deep Generative Models (CS236) Class

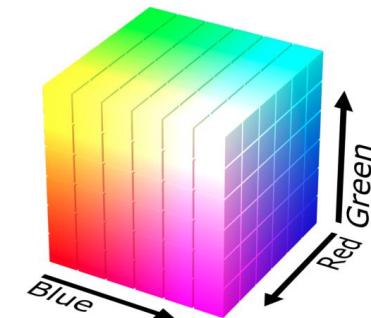
Sampling from the joint distribution randomly generates a color for the pixel:

$$r, g, b \sim p(R, G, B)$$

Image Example

Modeling a single pixel in an image \Rightarrow 3 discrete random variables:

- R - Red Channel. $\text{Val}(R) = \{0, \dots, 255\}$
- G - Green Channel. $\text{Val}(G) = \{0, \dots, 255\}$
- B - Blue Channel. $\text{Val}(B) = \{0, \dots, 255\}$



Source: Stefano Ermon, Deep Generative Models (CS236) Class

Sampling from the joint distribution randomly generates a color for the pixel:

$$r, g, b \sim p(R, G, B)$$

Recall from prob overview:
a sample from the joint PDF.

Image Example

Question: How many parameters do we need to specify the joint distribution

$$p(R = r, G = g, B = b) \quad ?$$

Image Example

Question: How many parameters do we need to specify the joint distribution

$$p(R = r, G = g, B = b) \quad ?$$

Answer: It's a discrete distribution with $256 \times 256 \times 256$ possible values \Rightarrow

$$256 \times 256 \times 256 - 1 = 16,777,215 \text{ parameters}$$

A lot!

Image Example

Modeling an image

Image Example

Modeling an image \Rightarrow For simplicity, consider black & white images with n pixels.

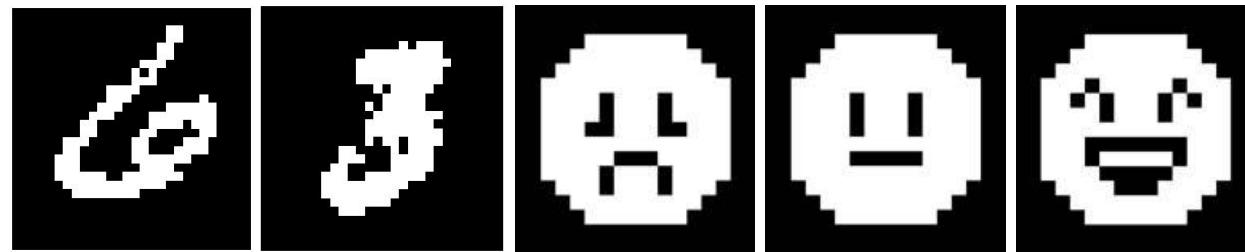
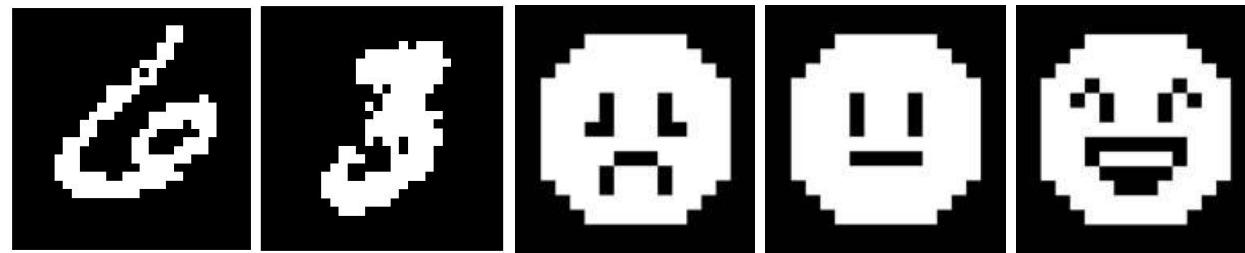


Image Example

Modeling an image \Rightarrow For simplicity, consider black & white images with n pixels.

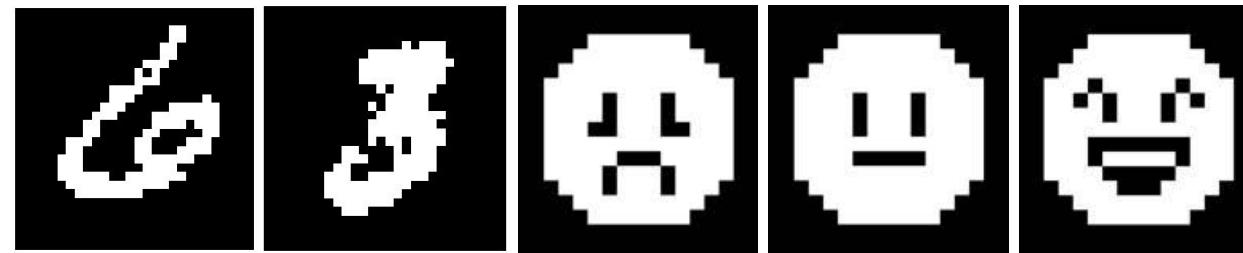


Model the image as a set of n random variables, X_1, X_2, \dots, X_n , where

$$\text{Val}(X_i) = \{0, 1\} = \{\text{Black}, \text{White}\}$$

Image Example

Modeling an image \Rightarrow For simplicity, consider black & white images with n pixels.



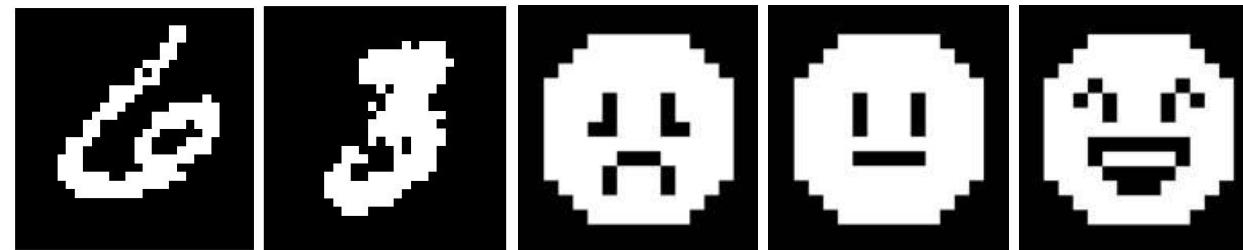
Model the image as a set of n random variables, X_1, X_2, \dots, X_n , where

$$\text{Val}(X_i) = \{0, 1\} = \{\text{Black}, \text{White}\}$$

How many possible images?

Image Example

Modeling an image \Rightarrow For simplicity, consider black & white images with n pixels.



Model the image as a set of n random variables, X_1, X_2, \dots, X_n , where

$$\text{Val}(X_i) = \{0, 1\} = \{\text{Black}, \text{White}\}$$

How many possible images? \Rightarrow $\underbrace{2 \times \dots \times 2}_{n \text{ times}} = 2^n$

Image Example

Question: How many parameters do we need to specify the joint distribution

$$p(X_1 = x_1, \dots, X_n = x_n) ?$$

Image Example

Question: How many parameters do we need to specify the joint distribution

$$p(X_1 = x_1, \dots, X_n = x_n) ?$$

Answer: It's a discrete distribution with 2^n possible values \Rightarrow

$$2^n - 1 \text{ parameters}$$

Image Example

Question: If we simplify our joint density, can we reduce this number?

Image Example

Question: What if we assume that X_1, X_2, \dots, X_n are independent?

Image Example

Question: What if we assume that X_1, X_2, \dots, X_n are independent?

Answer: Then $p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2) \cdots p(x_n)$.

Note: this is still a discrete distribution with 2^n possible values.

Image Example

Question: What if we assume that X_1, X_2, \dots, X_n are independent?

Answer: Then $p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2) \cdots p(x_n)$.

Note: this is still a discrete distribution with 2^n possible values.

However, to specify a marginal distribution $p(x_i)$ requires 1 parameter.

Binary random variable specified by 1 parameter.

Image Example

Question: What if we assume that X_1, X_2, \dots, X_n are independent?

Answer: Then $p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2) \cdots p(x_n)$.

Note: this is still a discrete distribution with 2^n possible values.

However, to specify a marginal distribution $p(x_i)$ requires 1 parameter.

⇒ Specifying $p(x_1, x_2, \dots, x_n)$ requires only n parameters.

Binary random variable specified by 1 parameter.

⇒ The frequency of 2^n values can be described using only n numbers

⇒ Big reduction!

Image Example

Question: What if we assume that X_1, X_2, \dots, X_n are independent?

However, independence might
be too strong of an assumption!

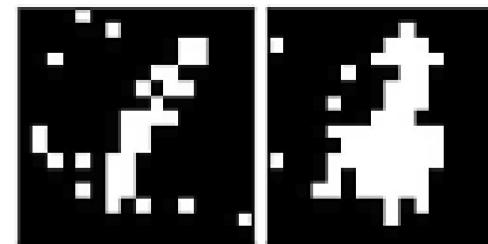
Image Example

Question: What if we assume that X_1, X_2, \dots, X_n are independent?

However, independence might be too strong of an assumption!

Model likely to be **not useful**.

For example, each pixel chosen independently when sampling does not yield realistic images.



vs.

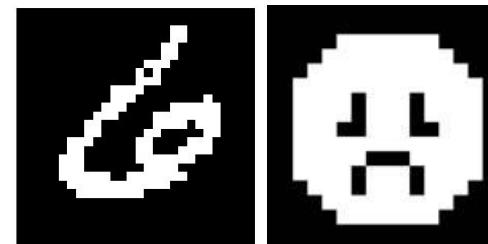


Image Example

Question: Can we think of another assumption that may be “in between”?

Image Example

Question: Can we think of another assumption that may be “in between”?

To answer this, remember a couple of properties from last class:

Image Example

Question: Can we think of another assumption that may be “in between”?

To answer this, remember a couple of properties from last class:

Joint PDF Factorization

Definition of conditional probability density function:

$$p(x_1 \mid x_2, \dots, x_n) = \frac{p(x_1, x_2, \dots, x_n)}{p(x_2, \dots, x_n)}$$

Image Example

Question: Can we think of another assumption that may be “in between”?

To answer this, remember a couple of properties from last class:

Joint PDF Factorization

Definition of conditional probability density function:

$$p(x_1 \mid x_2, \dots, x_n) = \frac{p(x_1, x_2, \dots, x_n)}{p(x_2, \dots, x_n)}$$

Leads to the **factorization** (i.e., way of writing) the joint PDF as:

$$p(x_1, x_2, \dots, x_n) = p(x_1 \mid x_2, \dots, x_n) p(x_2, \dots, x_n)$$

Image Example

Question: Can we think of another assumption that may be “in between”?

To answer this, remember a couple of properties from last class:

Chain Rule

Image Example

Question: Can we think of another assumption that may be “in between”?

To answer this, remember a couple of properties from last class:

Chain Rule

Joint PDF factorization

$$p(x_1, x_2, \dots, x_n) = p(x_n \mid x_1, x_2, \dots, x_{n-1})p(x_1, x_2, \dots, x_{n-1})$$

Image Example

Question: Can we think of another assumption that may be “in between”?

To answer this, remember a couple of properties from last class:

Chain Rule

Joint PDF factorization

$$\begin{aligned} p(x_1, x_2, \dots, x_n) &= p(x_n \mid x_1, x_2, \dots, x_{n-1}) p(x_1, x_2, \dots, x_{n-1}) && \text{(repeatedly!)} \\ &= p(x_n \mid x_1, x_2, \dots, x_{n-1}) p(x_{n-1} \mid x_1, x_2, \dots, x_{n-2}) p(x_1, x_2, \dots, x_{n-2}) \\ &= \dots \end{aligned}$$

Image Example

Question: Can we think of another assumption that may be “in between”?

To answer this, remember a couple of properties from last class:

Chain Rule

Joint PDF factorization

$$\begin{aligned} p(x_1, x_2, \dots, x_n) &= p(x_n \mid x_1, x_2, \dots, x_{n-1}) p(x_1, x_2, \dots, x_{n-1}) && \text{(repeatedly!)} \\ &= p(x_n \mid x_1, x_2, \dots, x_{n-1}) p(x_{n-1} \mid x_1, x_2, \dots, x_{n-2}) p(x_1, x_2, \dots, x_{n-2}) \\ &= \dots \\ &= p(x_1) \prod_{i=2}^n p(x_i \mid x_1, \dots, x_{i-1}) \\ &= p(x_1) p(x_2 \mid x_1) p(x_3 \mid x_2, x_1) \cdots p(x_n \mid x_{n-1}, \dots, x_1) \end{aligned}$$

Image Example

Question: Can we think of another assumption that may be “in between”?

Image Example

Question: Can we think of another assumption that may be “in between”?

Answer:

⇒ First, Consider the chain rule:

$$p(x_1, x_2, \dots, x_n) = p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1) \cdots p(x_n | x_{n-1}, \dots, x_1)$$

Image Example

Question: Can we think of another assumption that may be “in between”?

Answer:

$p(x_1)$ requires 1 parameter

⇒ First, Consider the chain rule:

$$p(x_1, x_2, \dots, x_n) = p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1) \cdots p(x_n | x_{n-1}, \dots, x_1)$$

How many parameters? $1 + 2 + \dots + 2^{n-1} = 2^n - 1$

Image Example

Question: Can we think of another assumption that may be “in between”?

Answer:

$p(x_2 \mid x_1 = 0)$ requires 1 parameter
 $p(x_2 \mid x_1 = 1)$ requires 1 parameter

⇒ First, Consider the chain rule:

$$p(x_1, x_2, \dots, x_n) = p(x_1) p(x_2 \mid x_1) p(x_3 \mid x_2, x_1) \cdots p(x_n \mid x_{n-1}, \dots, x_1)$$

How many parameters? $1 + 2 + \dots + 2^{n-1} = 2^n - 1$

Image Example

Question: Can we think of another assumption that may be “in between”?

Answer:

⇒ First, Consider the chain rule:

$$p(x_1, x_2, \dots, x_n) = p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1) \cdots p(x_n | x_{n-1}, \dots, x_1)$$

How many parameters? $1 + 2 + \dots + 2^{n-1} = 2^n - 1$

(So this is the same as before, when we had no independence assumptions!)

⇒ the chain rule does not buy us anything.

Image Example

Question: Can we think of another assumption that may be “in between”?

Answer:

Now suppose each X_i is independent of variables $X_1 \dots X_{i-2}$ given X_{i-1} , i.e.,

$$X_i \perp (X_1 \dots X_{i-2}) \mid X_{i-1}$$

Image Example

Question: Can we think of another assumption that may be “in between”?

Answer:

Now suppose each X_i is independent of variables $X_1 \dots X_{i-2}$ given X_{i-1} , i.e.,

$$X_i \perp (X_1 \dots X_{i-2}) \mid X_{i-1}$$

Could call this a “Markovian” dependence structure.

Image Example

Question: Can we think of another assumption that may be “in between”?

Answer:

Now suppose each X_i is independent of variables $X_1 \dots X_{i-2}$ given X_{i-1} , i.e.,

$$X_i \perp (X_1 \dots X_{i-2}) \mid X_{i-1}$$

Then:

$$\begin{aligned} p(x_1, x_2, \dots, x_n) &= p(x_1) p(x_2 \mid x_1) p(x_3 \mid x_2, x_1) \cdots p(x_n \mid x_{n-1}, \dots, x_1) \\ &= p(x_1) p(x_2 \mid x_1) p(x_3 \mid x_2, \cancel{x_1}) \cdots p(x_n \mid x_{n-1}, \cancel{\dots}, \cancel{x_1}) \\ &= p(x_1) p(x_2 \mid x_1) p(x_3 \mid x_2) \cdots p(x_n \mid x_{n-1}) \end{aligned}$$

Image Example

Question: Can we think of another assumption that may be “in between”?

Answer:

Now suppose each X_i is independent of variables $X_1 \dots X_{i-2}$ given X_{i-1} , i.e.,

$$X_i \perp (X_1 \dots X_{i-2}) \mid X_{i-1}$$

Then:

$$p(x_1, x_2, \dots, x_n) = p(x_1) p(x_2 \mid x_1) p(x_3 \mid x_2) \cdots p(x_n \mid x_{n-1})$$

⇒ In total, $2n - 1$ parameters.

⇒ Exponential reduction!

Conditional Independence Assumptions

We can incorporate structure (simplifications) through conditional independence!
i.e., conditional independence assumptions

Conditional Independence Assumptions

We can incorporate structure (simplifications) through conditional independence!
i.e., conditional independence assumptions

Bayesian Networks – the general idea:

- Represent a joint PDF (probability model) in terms of a set of conditional parameterizations.
- For each random variable X_i , specify a distribution $p(x_i \mid \mathbf{x}_{P_i})$ for a set of variables \mathbf{x}_{P_i} .
- Then write the joint parameterization as:

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i \mid \mathbf{x}_{P_i})$$

(Need to guarantee this is a “legal” probability distribution → must correspond to some chain rule factorization with factors simplified due to assumed conditional independencies).

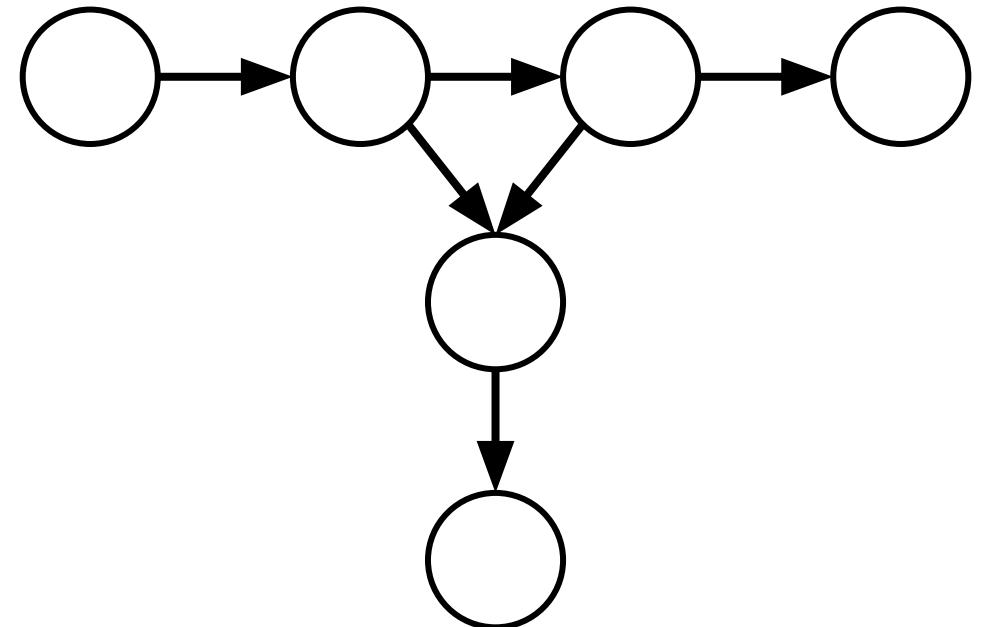
Bayesian Networks

Definition.

Bayesian Networks

Definition. A Bayesian Network is defined by a directed acyclic graph (DAG)
Where (V, E)

1. There is one node $i \in V$ for each random variable X_i .
2. There is one conditional probability distribution per node, $p(x_i | \mathbf{x}_{\text{Pa}(i)})$, specifying the variable's probability conditioned on its parents' values.



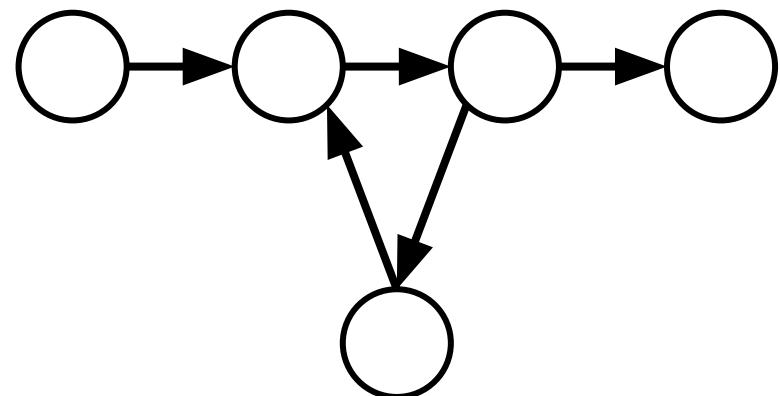
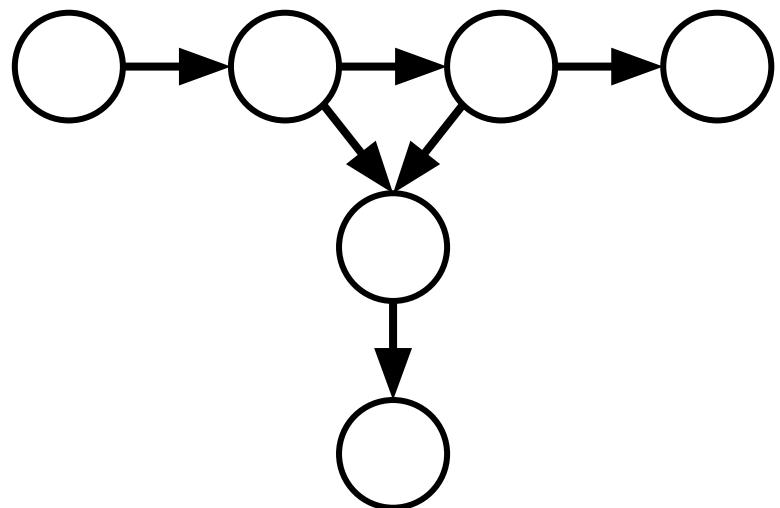
Bayesian Networks

Note.

Directed acyclic graph (DAG)

vs.

Not-a-DAG.



Bayesian Networks

The graph $G = (V, E)$ is called the structure of the Bayesian network.

This defines a joint distribution:

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$$

Bayesian Networks

The graph $G = (V, E)$ is called the structure of the Bayesian network.

This defines a joint distribution:

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$$

Claim:

$p(x_1, \dots, x_n)$ is a valid probability distribution (simplifies some chain rule rep'n), because of the ordering implied by the DAG.

Bayesian Networks

The graph $G = (V, E)$ is called the structure of the Bayesian network.

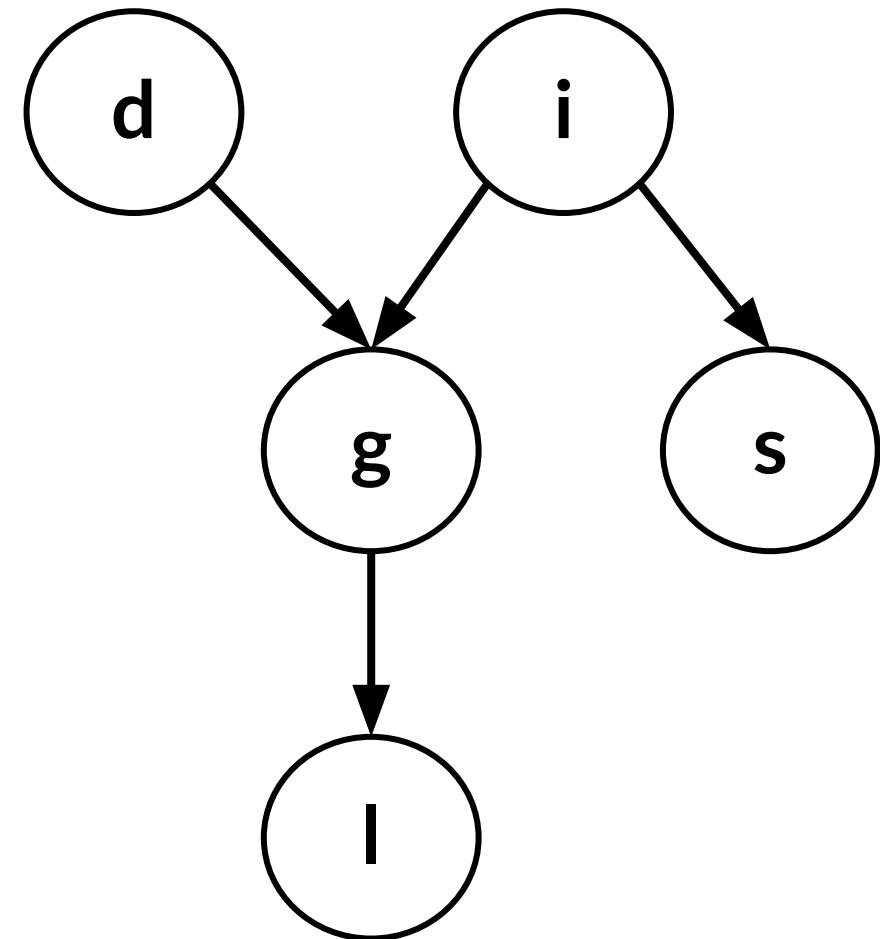
This defines a joint distribution:

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$$

Economical representation \Rightarrow Number of parameters is exponential in $\max_{i \in V} |\text{Pa}(i)|$, not in $|V|$.

Bayesian Networks – Example

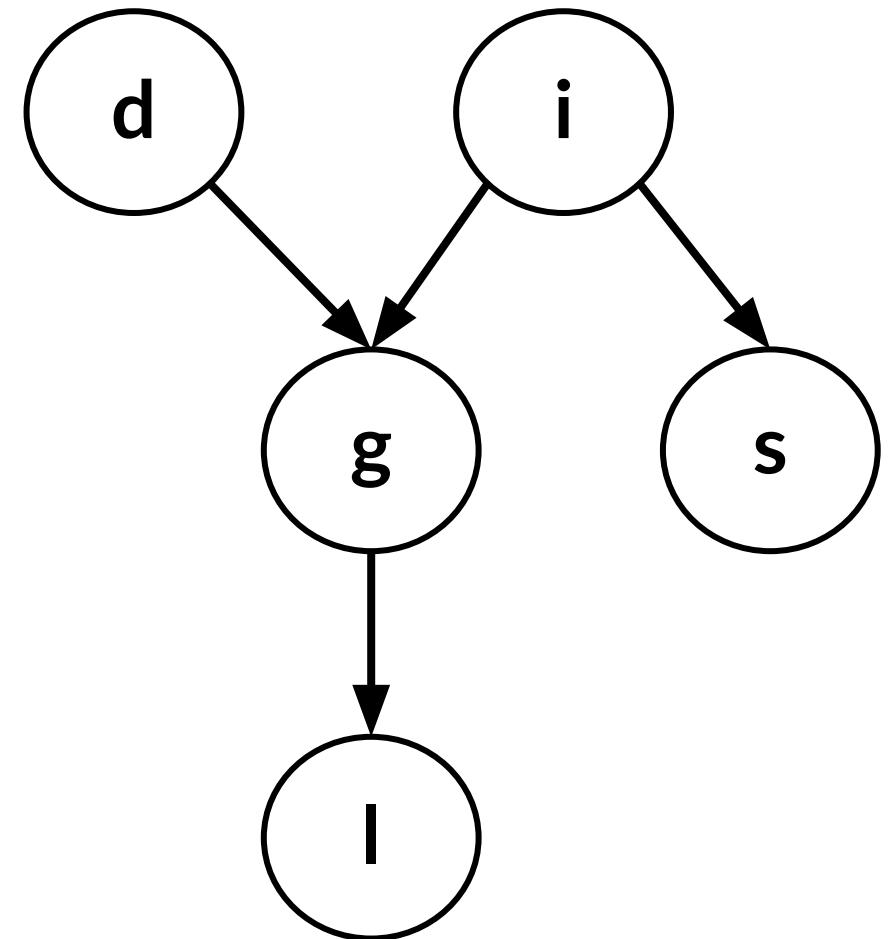
Consider the following Bayesian network:



Bayesian Networks – Example

Consider the following Bayesian network:

(A model of student grades and recommendation letters)

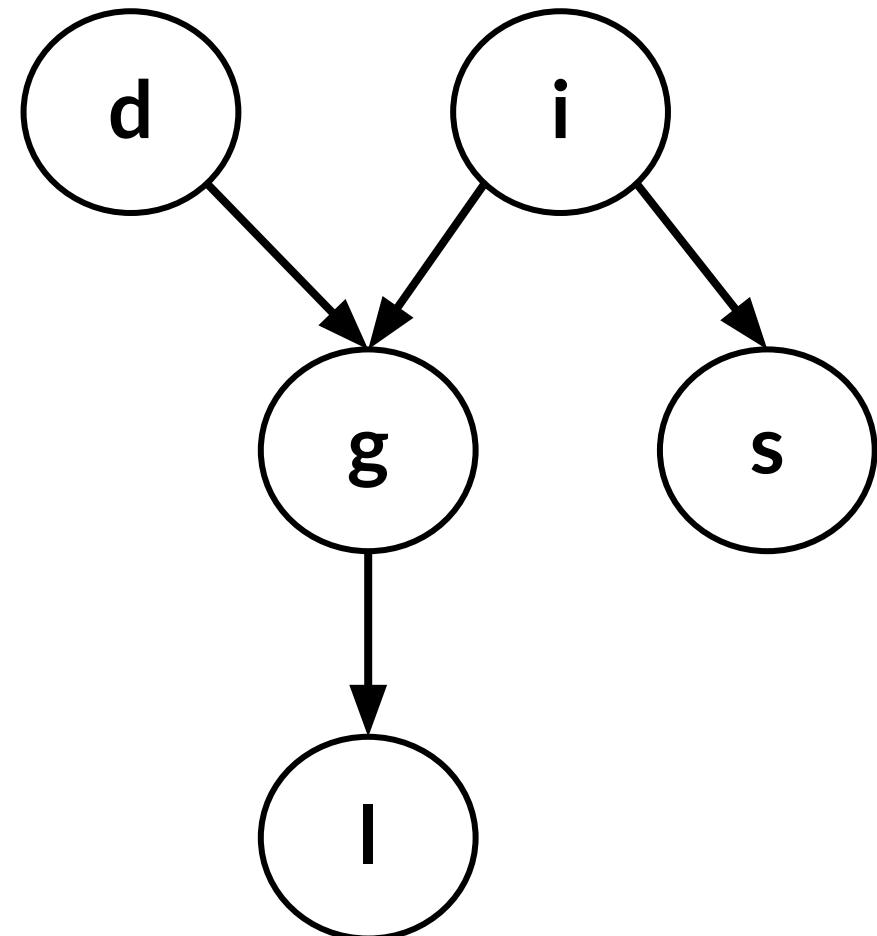


Bayesian Networks – Example

Consider the following Bayesian network:

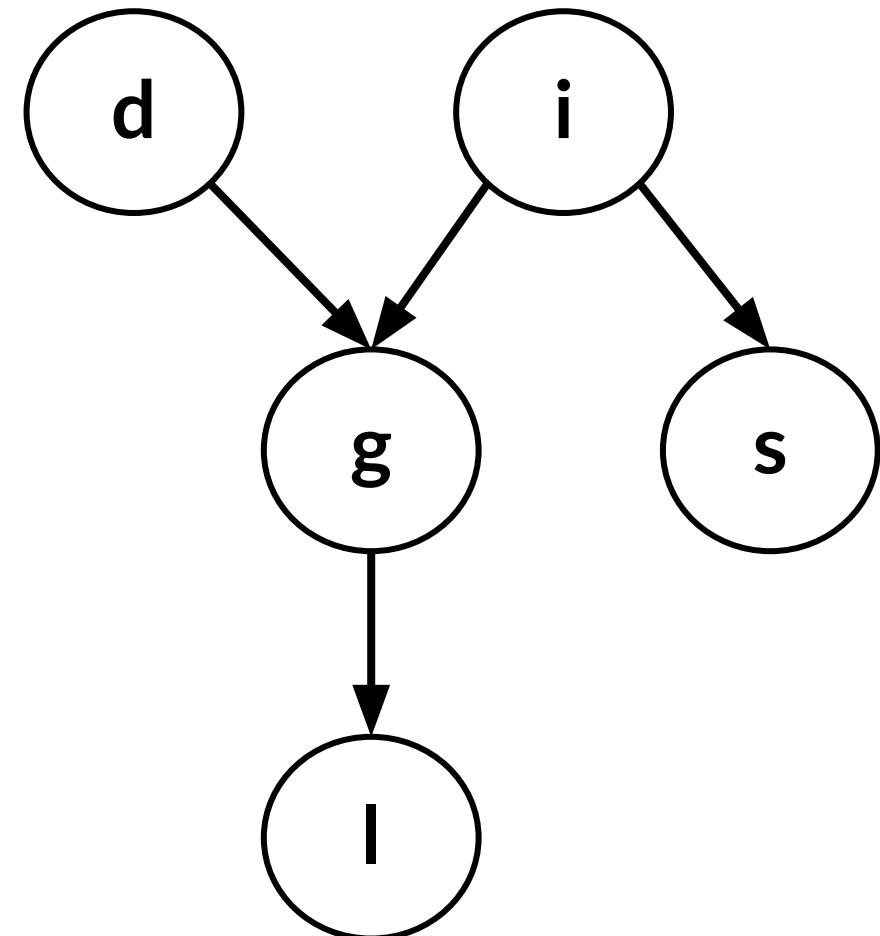
(A model of student grades and recommendation letters)

- d - Class difficulty.
- i - Student's intelligence.
- g - Student's grade in the class.
- s - Student's SAT test score.
- l - Professor's letter of recommendation.
(good vs bad letter)



Bayesian Networks – Example

Think of it as a generative process that generates data from top to bottom
(i.e., from parents to children)



Bayesian Networks – Example

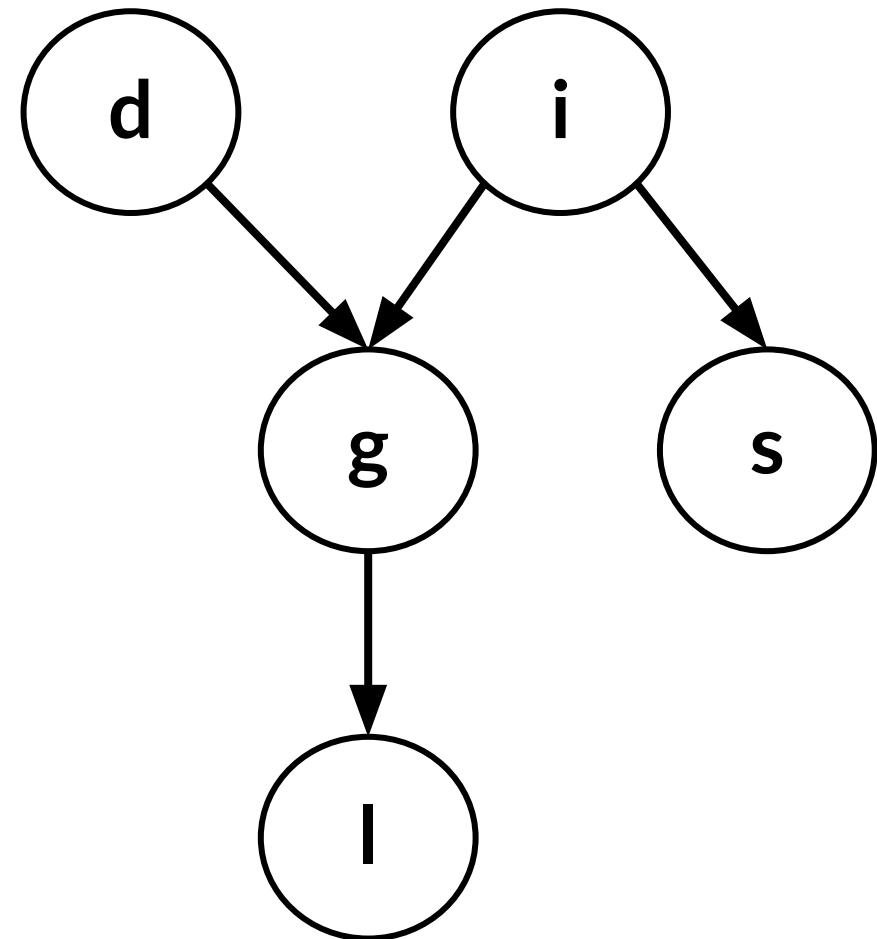
Think of it as a generative process that generates data from top to bottom (i.e., from parents to children)

Namely:

Class difficulty d and student's intelligence i → student's grade g .

Student's grade g → letter of recommendation I .

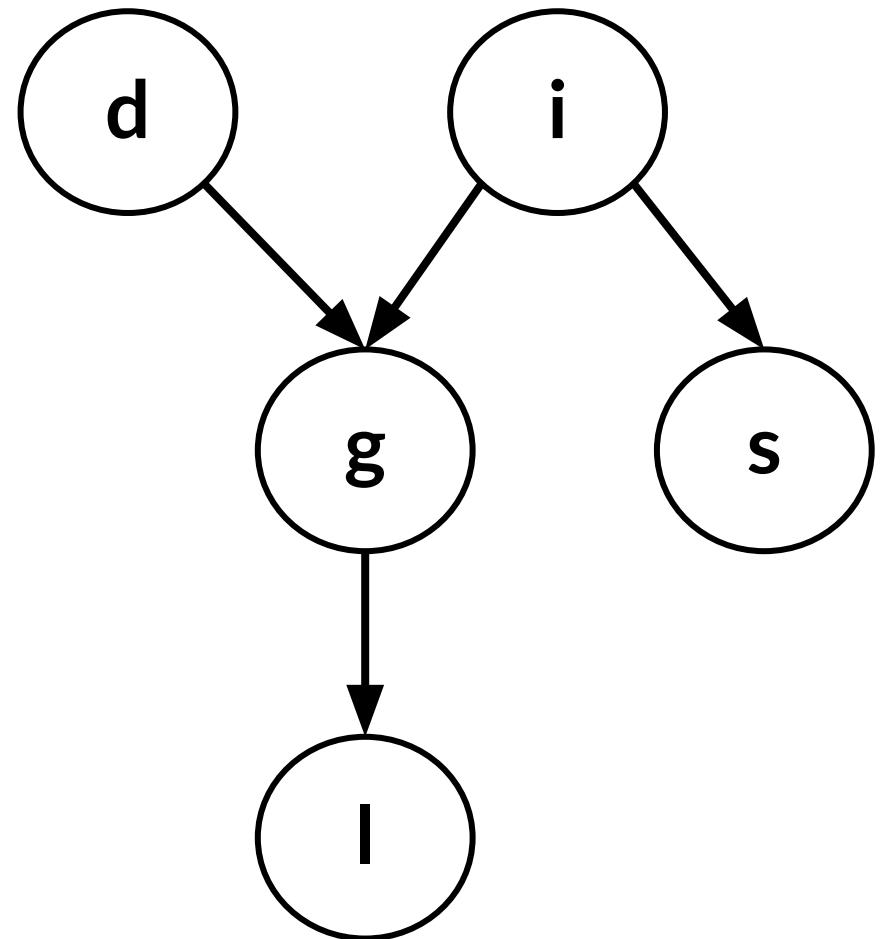
Student's intelligence → SAT score s



Bayesian Networks – Example

How can we parameterize
these distributions?

*I.e., “define”
or “specify”*



Bayesian Networks – Example

How can we parameterize these distributions?

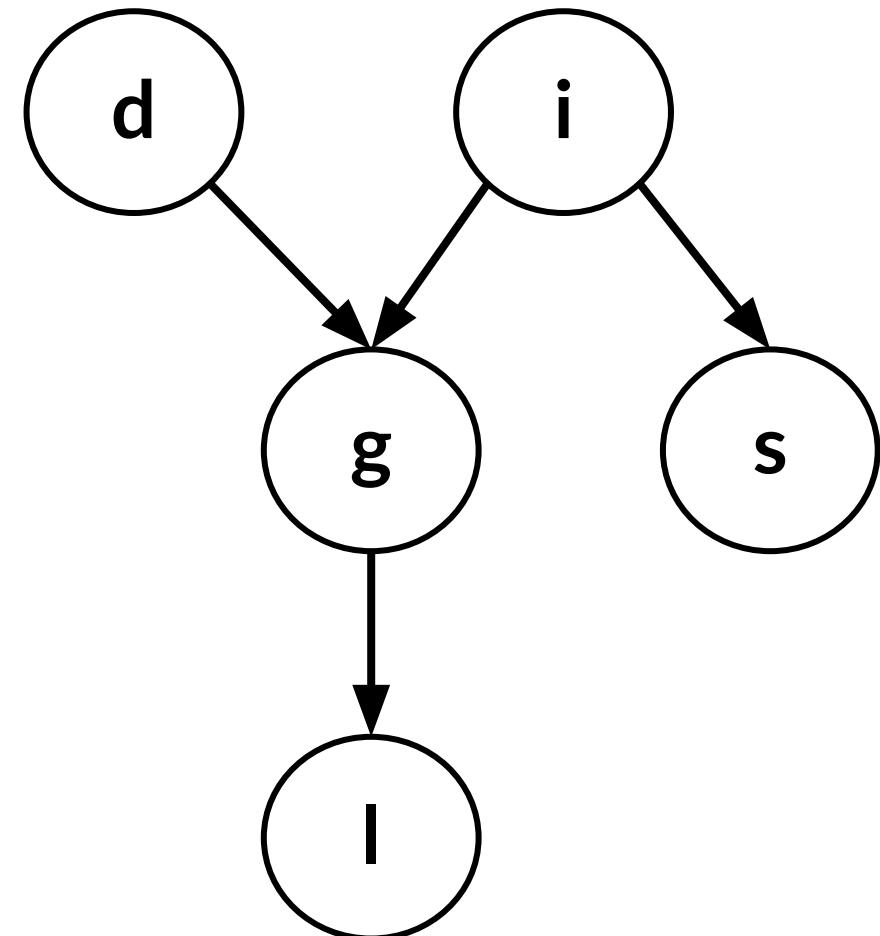
Suppose

d, i, s, l – binary random variables

- d - easy vs hard class.
- i - intelligent vs not-as-intelligent.
- s - good SAT score vs bad SAT score.
- l - good letter of rec vs bad letter of rec.

g – has three values.

- g - good grade vs average grade vs bad grade.



Bayesian Networks – Example

Can define conditional probability distributions!

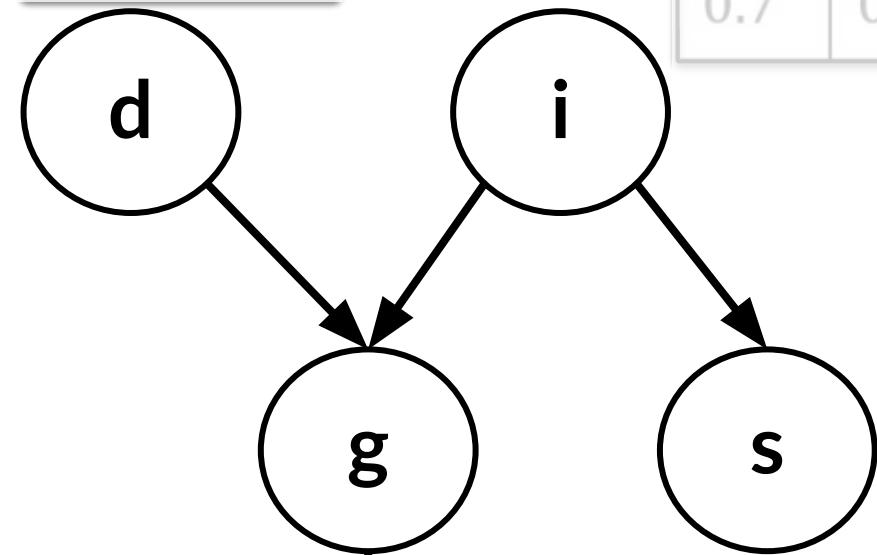
(In this case – tables).

| | g^1 | g^2 | g^3 |
|------------|-------|-------|-------|
| i^0, d^0 | 0.3 | 0.4 | 0.3 |
| i^0, d^1 | 0.05 | 0.25 | 0.7 |
| i^1, d^0 | 0.9 | 0.08 | 0.02 |
| i^1, d^1 | 0.5 | 0.3 | 0.2 |

| | l^0 | l^1 |
|-------|-------|-------|
| g^1 | 0.1 | 0.9 |
| g^2 | 0.4 | 0.6 |
| g^3 | 0.99 | 0.01 |

| d^0 | d^1 |
|-------|-------|
| 0.6 | 0.4 |

| i^0 | i^1 |
|-------|-------|
| 0.7 | 0.3 |



| | s^0 | s^1 |
|-------|-------|-------|
| i^0 | 0.95 | 0.05 |
| i^1 | 0.2 | 0.8 |

Bayesian Networks – Example

Can define conditional probability distributions!

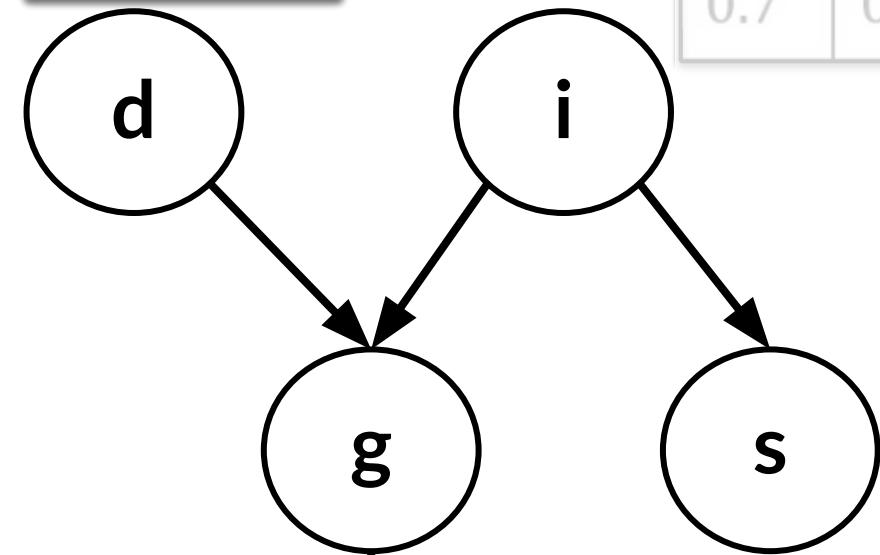
(In this case – tables).

| | g^1 | g^2 | g^3 |
|------------|-------|-------|-------|
| i^0, d^0 | 0.3 | 0.4 | 0.3 |
| i^0, d^1 | 0.05 | 0.25 | 0.7 |
| i^1, d^0 | 0.9 | 0.08 | 0.02 |
| i^1, d^1 | 0.5 | 0.3 | 0.2 |

| | l^0 | l^1 |
|-------|-------|-------|
| g^1 | 0.1 | 0.9 |
| g^2 | 0.4 | 0.6 |
| g^3 | 0.99 | 0.01 |

| d^0 | d^1 |
|-------|-------|
| 0.6 | 0.4 |

| i^0 | i^1 |
|-------|-------|
| 0.7 | 0.3 |



| | s^0 | s^1 |
|-------|-------|-------|
| i^0 | 0.95 | 0.05 |
| i^1 | 0.2 | 0.8 |

Bayesian Networks – Example

Can define conditional probability distributions!

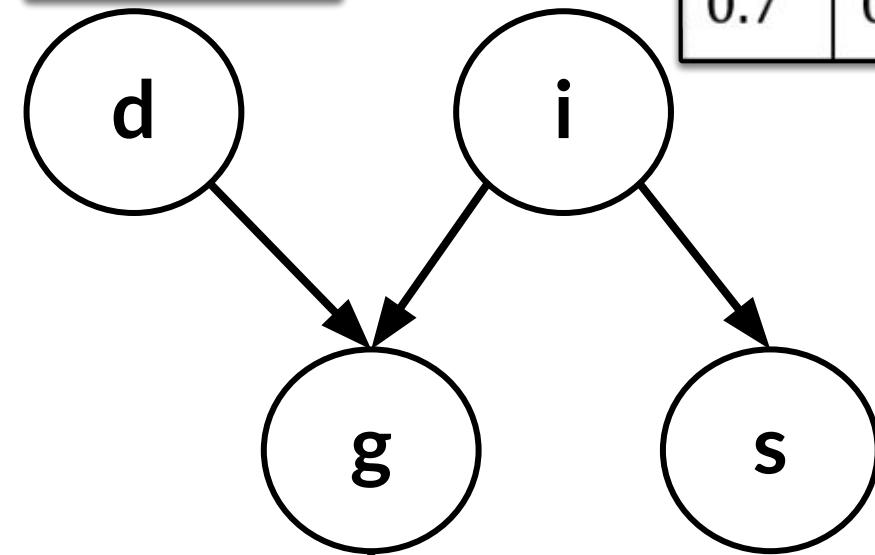
(In this case – tables).

| | g^1 | g^2 | g^3 |
|------------|-------|-------|-------|
| i^0, d^0 | 0.3 | 0.4 | 0.3 |
| i^0, d^1 | 0.05 | 0.25 | 0.7 |
| i^1, d^0 | 0.9 | 0.08 | 0.02 |
| i^1, d^1 | 0.5 | 0.3 | 0.2 |

| | l^0 | l^1 |
|-------|-------|-------|
| g^1 | 0.1 | 0.9 |
| g^2 | 0.4 | 0.6 |
| g^3 | 0.99 | 0.01 |

| d^0 | d^1 |
|-------|-------|
| 0.6 | 0.4 |

| i^0 | i^1 |
|-------|-------|
| 0.7 | 0.3 |



| | s^0 | s^1 |
|-------|-------|-------|
| i^0 | 0.95 | 0.05 |
| i^1 | 0.2 | 0.8 |

Bayesian Networks – Example

Can define conditional probability distributions!

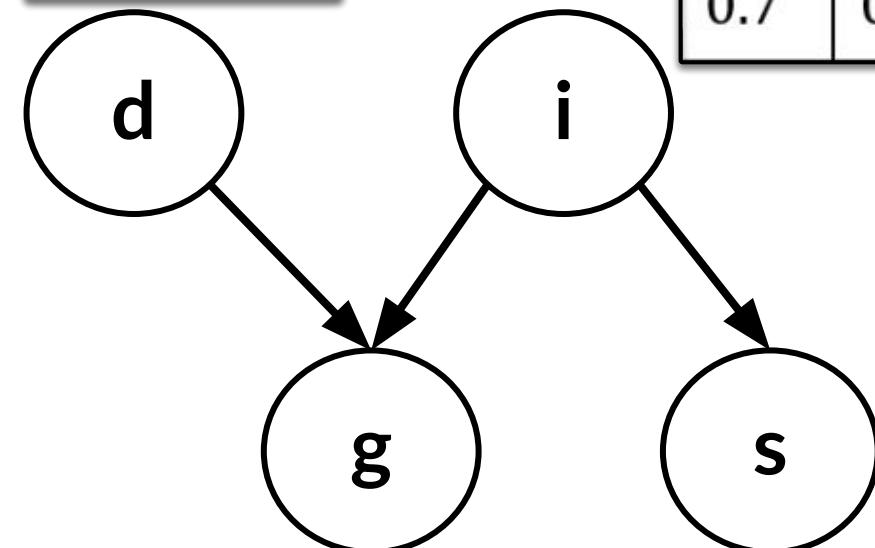
(In this case – tables).

| | g^1 | g^2 | g^3 |
|------------|-------|-------|-------|
| i^0, d^0 | 0.3 | 0.4 | 0.3 |
| i^0, d^1 | 0.05 | 0.25 | 0.7 |
| i^1, d^0 | 0.9 | 0.08 | 0.02 |
| i^1, d^1 | 0.5 | 0.3 | 0.2 |

| | l^0 | l^1 |
|-------|-------|-------|
| g^1 | 0.1 | 0.9 |
| g^2 | 0.4 | 0.6 |
| g^3 | 0.99 | 0.01 |

| d^0 | d^1 |
|-------|-------|
| 0.6 | 0.4 |

| i^0 | i^1 |
|-------|-------|
| 0.7 | 0.3 |



| | s^0 | s^1 |
|-------|-------|-------|
| i^0 | 0.95 | 0.05 |
| i^1 | 0.2 | 0.8 |

Bayesian Networks – Example

Can define conditional probability distributions!

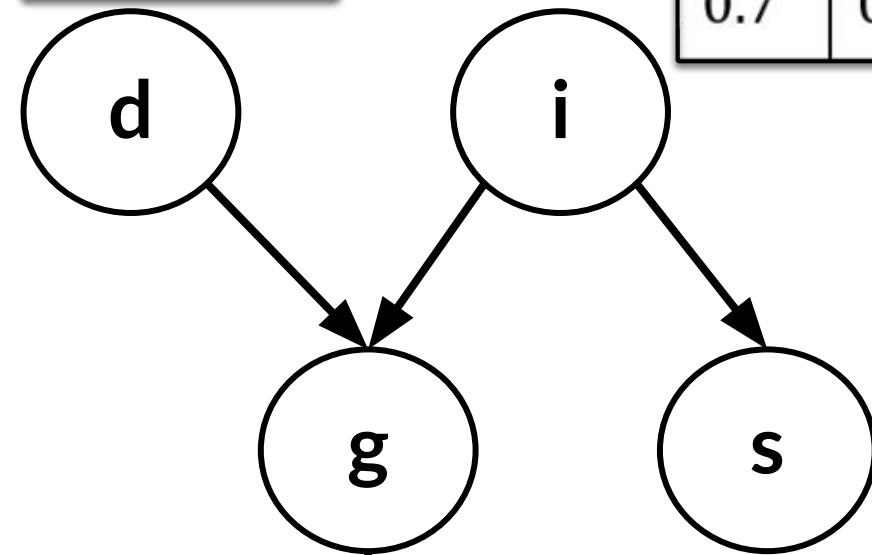
(In this case – tables).

| | g^1 | g^2 | g^3 |
|------------|-------|-------|-------|
| i^0, d^0 | 0.3 | 0.4 | 0.3 |
| i^0, d^1 | 0.05 | 0.25 | 0.7 |
| i^1, d^0 | 0.9 | 0.08 | 0.02 |
| i^1, d^1 | 0.5 | 0.3 | 0.2 |

| | l^0 | l^1 |
|-------|-------|-------|
| g^1 | 0.1 | 0.9 |
| g^2 | 0.4 | 0.6 |
| g^3 | 0.99 | 0.01 |

| d^0 | d^1 |
|-------|-------|
| 0.6 | 0.4 |

| i^0 | i^1 |
|-------|-------|
| 0.7 | 0.3 |



| | s^0 | s^1 |
|-------|-------|-------|
| i^0 | 0.95 | 0.05 |
| i^1 | 0.2 | 0.8 |

Bayesian Networks – Example

Can define conditional probability distributions!

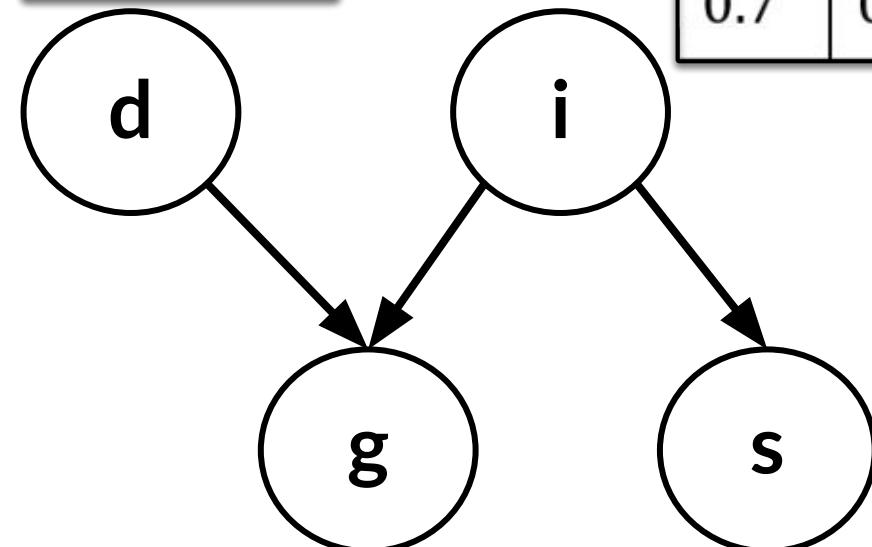
(In this case – tables).

| | g^1 | g^2 | g^3 |
|------------|-------|-------|-------|
| i^0, d^0 | 0.3 | 0.4 | 0.3 |
| i^0, d^1 | 0.05 | 0.25 | 0.7 |
| i^1, d^0 | 0.9 | 0.08 | 0.02 |
| i^1, d^1 | 0.5 | 0.3 | 0.2 |

| | l^0 | l^1 |
|-------|-------|-------|
| g^1 | 0.1 | 0.9 |
| g^2 | 0.4 | 0.6 |
| g^3 | 0.99 | 0.01 |

| d^0 | d^1 |
|-------|-------|
| 0.6 | 0.4 |

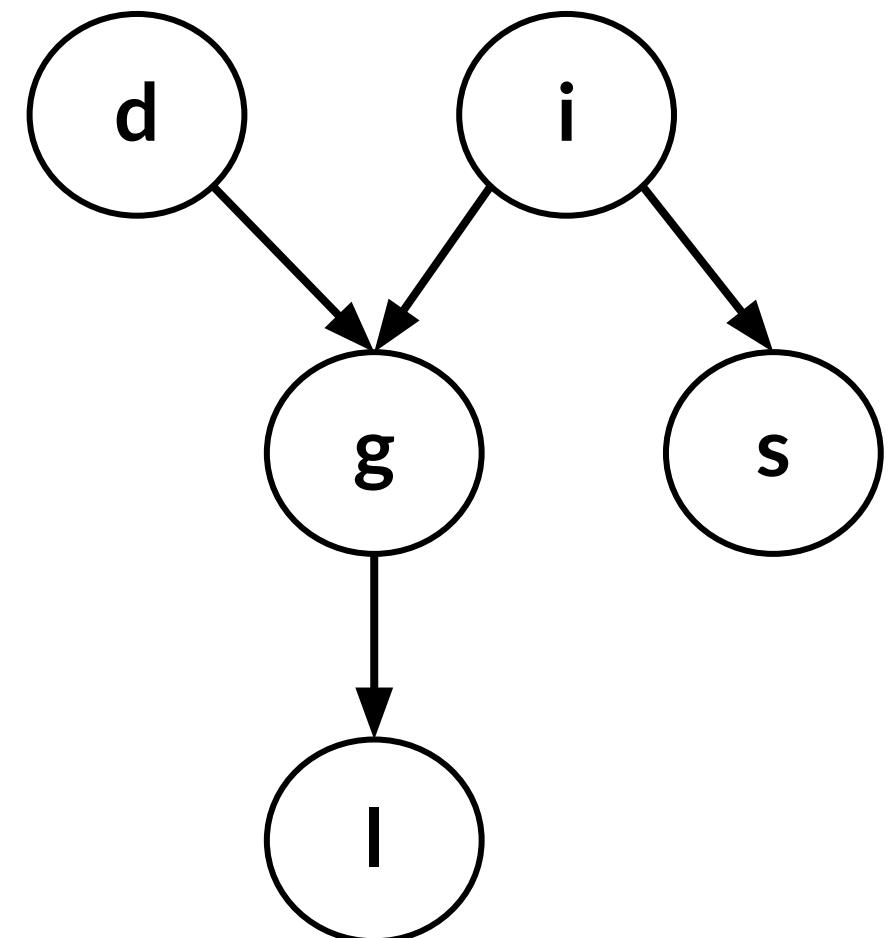
| i^0 | i^1 |
|-------|-------|
| 0.7 | 0.3 |



| | s^0 | s^1 |
|-------|-------|-------|
| i^0 | 0.95 | 0.05 |
| i^1 | 0.2 | 0.8 |

Bayesian Networks – Example

What is the joint distribution?

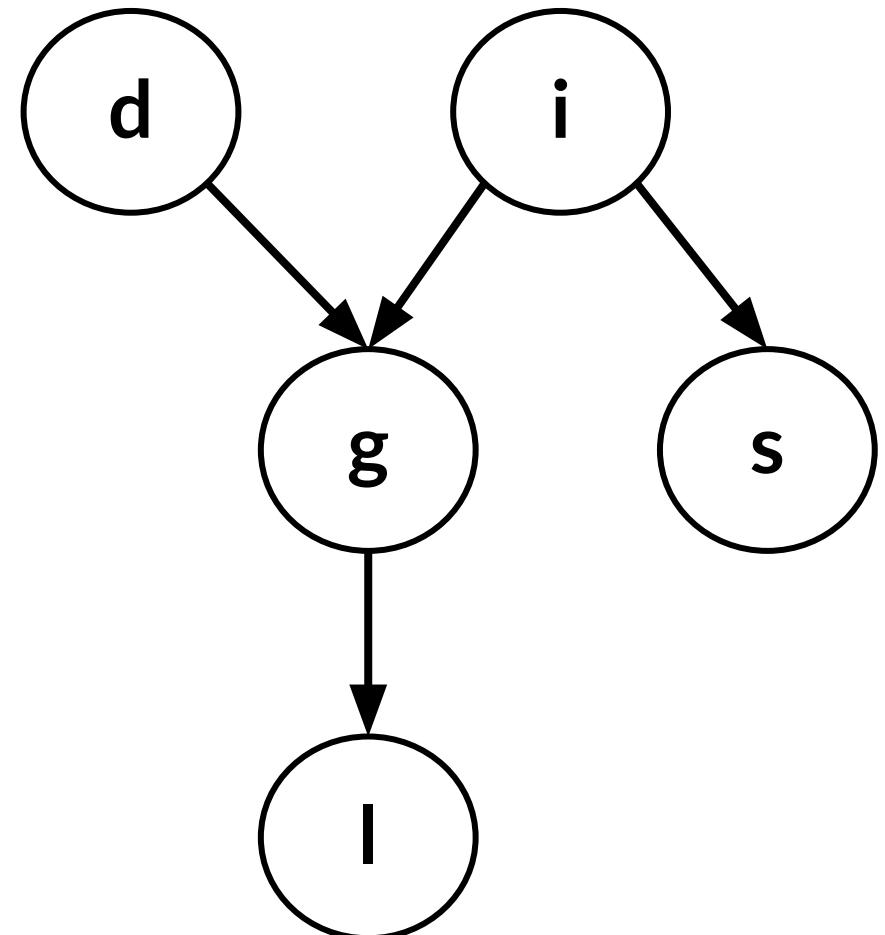


Bayesian Networks – Example

What is the joint distribution?

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i | \mathbf{x}_{\text{Pa}(i)})$$

$$p(l, g, i, d, s) = p(l | g) p(g | i, d) p(i) p(d) p(s | i)$$



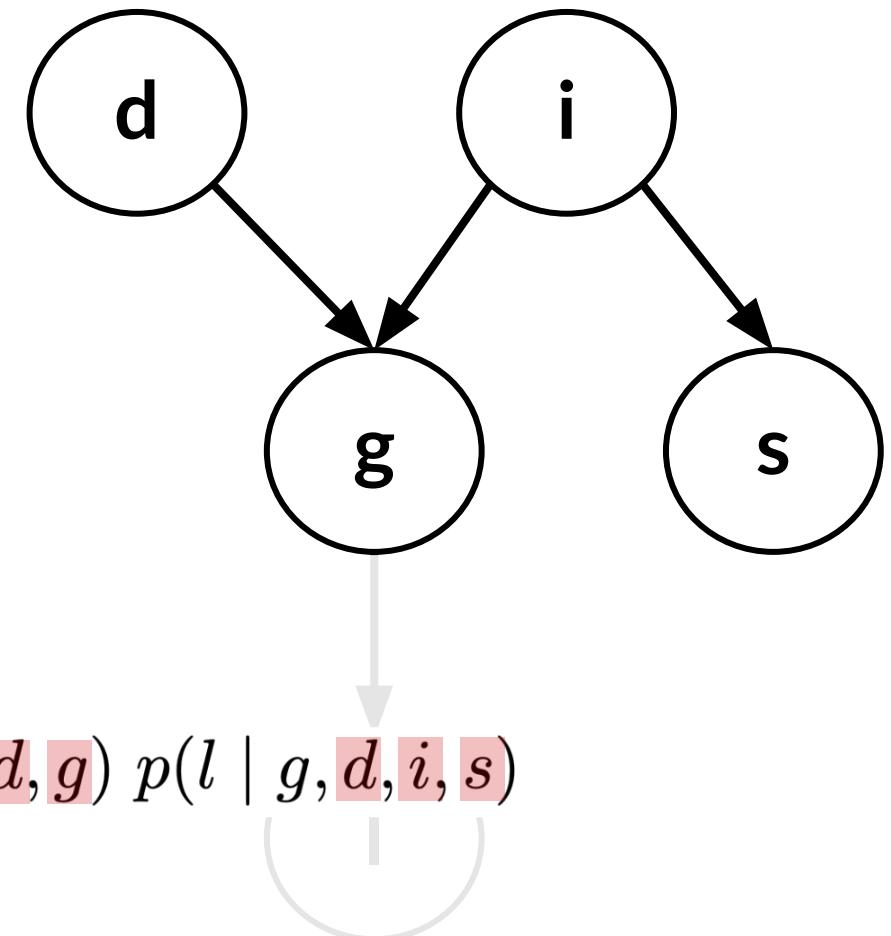
Bayesian Networks – Example

Joint distribution (re-written):

$$p(d, i, g, s, l) = p(d) p(i) p(g | i, d) p(s | i) p(l | g)$$

By chain rule (without any assumptions), could write joint distribution as:

$$p(d, i, g, s, l) = p(d) p(i | d) p(g | i, d) p(s | i, d, g) p(l | g, d, i, s)$$



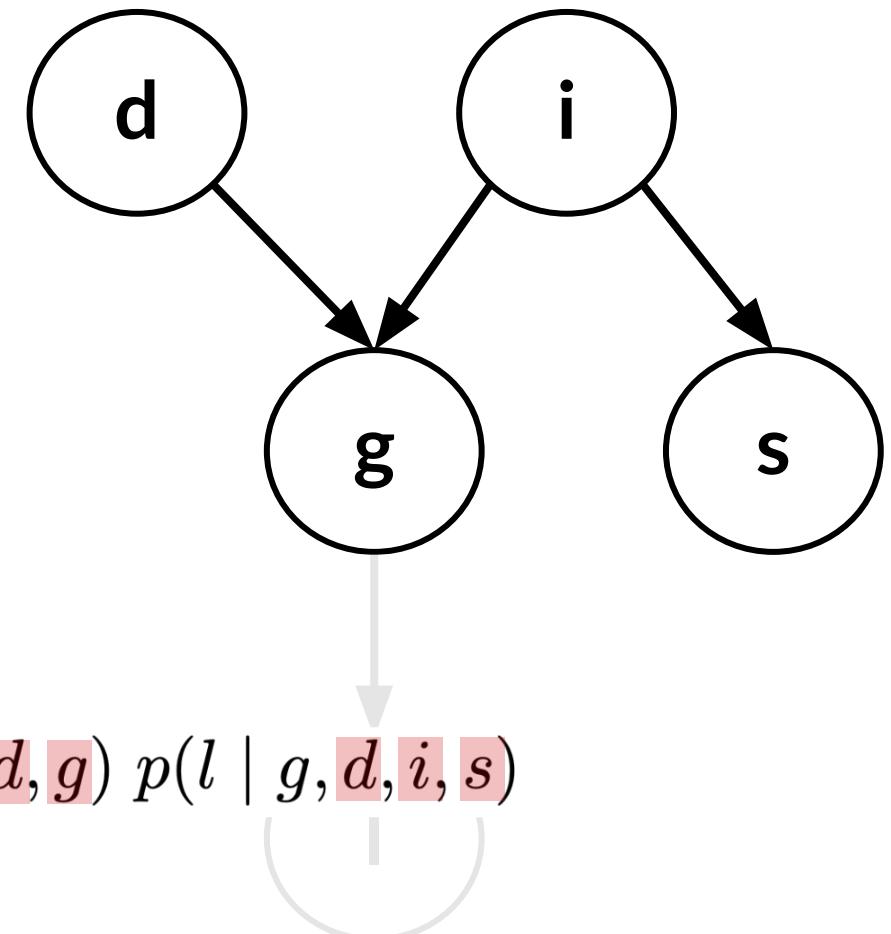
Bayesian Networks – Example

Joint distribution (re-written):

$$p(d, i, g, s, l) = p(d) p(i) p(g | i, d) p(s | i) p(l | g)$$

By chain rule (without any assumptions), could write joint distribution as:

$$p(d, i, g, s, l) = p(d) p(i | d) p(g | i, d) p(s | i, d, g) p(l | g, d, i, s)$$

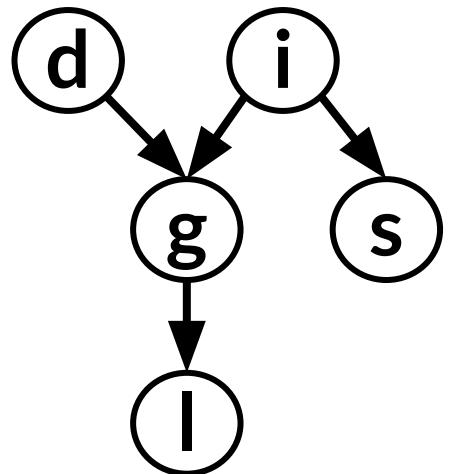


Thus, we are assuming the following conditional independencies:

$$d \perp i \quad s \perp \{d, g\} | i \quad l \perp \{i, d, s\} | g$$

Bayesian Networks – Generative Process View

Generative process view.

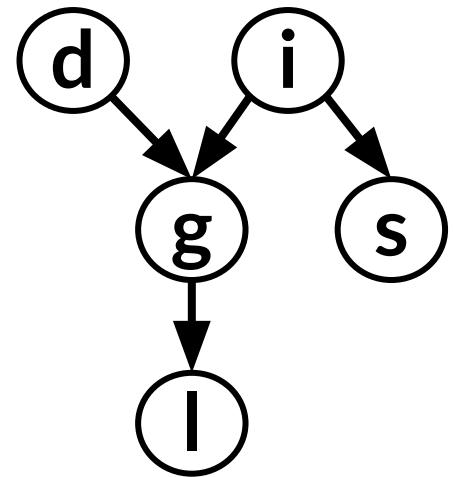


Bayesian Networks – Generative Process View

Generative process view.

View a Bayes net as “generating” from root to leaf variables.

- *I.e.,* view it as defining a *sampling procedure* for all variables in the graph.



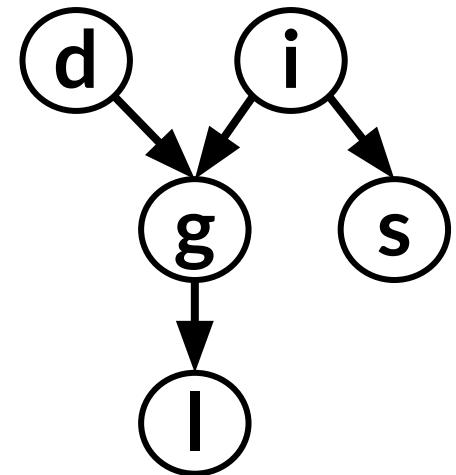
Bayesian Networks – Generative Process View

Generative process view.

View a Bayes net as “generating” from root to leaf variables.

- *I.e.,* view it as defining a *sampling procedure* for all variables in the graph.

Iteratively sample from conditional probability distribution according to a topological sort of the nodes.



Topological Sort (of a directed graph):

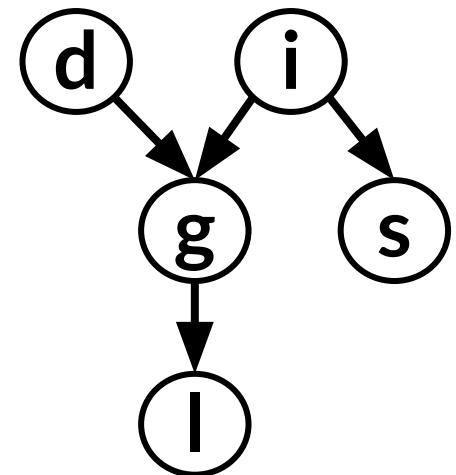
A linear ordering of the vertices such that for every directed edge (u, v) from vertex u to vertex v , u comes before v in the ordering

Bayesian Networks – Generative Process View

Generative process view.

View a Bayes net as “generating” from root to leaf variables.

- *I.e.*, view it as defining a *sampling procedure* for all variables in the graph.



Iteratively sample from conditional probability distribution according to a topological sort of the nodes.

$$\begin{array}{c} d \sim p(d) \qquad \qquad i \sim p(i) \\ \qquad \qquad \qquad \downarrow \qquad \qquad \qquad \downarrow \\ g \sim p(g \mid i, d) \qquad \qquad s \sim p(s \mid i) \\ \qquad \qquad \qquad \downarrow \\ l \sim p(l \mid g) \end{array}$$

Bayesian Networks – Generative Process View

Generative process view.

View a Bayes net as “generating” from root to leaf variables.

- *I.e.*, view it as defining a *sampling procedure* for all variables in the graph.

Iteratively sample from conditional probability distribution according to a topological sort of the nodes.

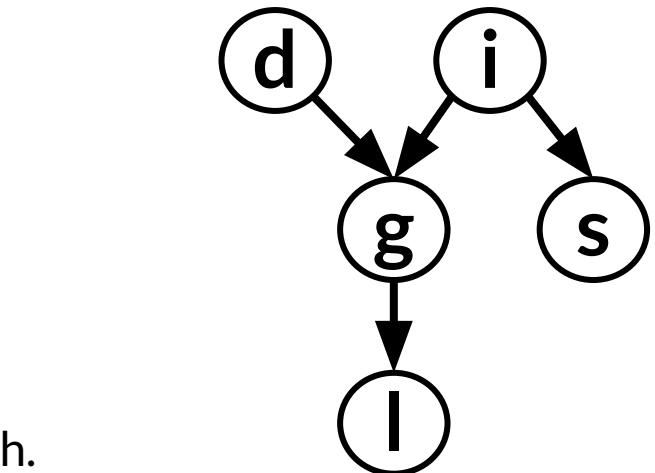
$$d \sim p(d)$$

$$i \sim p(i)$$

$$g \sim p(g \mid i, d)$$

$$s \sim p(s \mid i)$$

$$l \sim p(l \mid g)$$



Often called the
“**generative process**” view
of the PGM/joint PDF.

Bayesian Networks – Summary

Bayesian Networks – Summary

Bayesian network is defined by a tuple $(G = \{V, E\}, P)$ where P is specified as a set of local conditional probability distributions (CPDs) associated with nodes V

⇒ Efficient representation using a graph-based data structure.

Bayesian Networks – Summary

Bayesian network is defined by a tuple $(G = \{V, E\}, P)$ where P is specified as a set of local conditional probability distributions (CPDs) associated with nodes V

⇒ Efficient representation using a graph-based data structure.

Computing joint probability of any assignment obtained by multiplying CPDs.

Bayesian Networks – Summary

Bayesian network is defined by a tuple $(G = \{V, E\}, P)$ where P is specified as a set of local conditional probability distributions (CPDs) associated with nodes V

⇒ Efficient representation using a graph-based data structure.

Computing joint probability of any assignment obtained by multiplying CPDs.

Can sample from joint by sampling from CPDs according to DAG ordering.

Bayesian Networks – Summary

Bayesian network is defined by a tuple $(G = \{V, E\}, P)$ where P is specified as a set of local conditional probability distributions (CPDs) associated with nodes V

⇒ Efficient representation using a graph-based data structure.

Computing joint probability of any assignment obtained by multiplying CPDs.

Can sample from joint by sampling from CPDs according to DAG ordering.

Can identify some conditional independencies by looking at graph properties.

Bayesian Networks – Summary

Bayesian network is defined by a tuple $(G = \{V, E\}, P)$ where P is specified as a set of local conditional probability distributions (CPDs) associated with nodes V

⇒ Efficient representation using a graph-based data structure.

Computing joint probability of any assignment obtained by multiplying CPDs.

Can sample from joint by sampling from CPDs according to DAG ordering.

Can identify some conditional independencies by looking at graph properties.

Nodes in graph can correspond to random variables or random vectors.

Bayesian Networks – Summary

Bayesian network is defined by a tuple $(G = \{V, E\}, P)$ where P is specified as a set of local conditional probability distributions (CPDs) associated with nodes V

⇒ Efficient representation using a graph-based data structure.

Computing joint probability of any assignment obtained by multiplying CPDs.

Can sample from joint by sampling from CPDs according to DAG ordering.

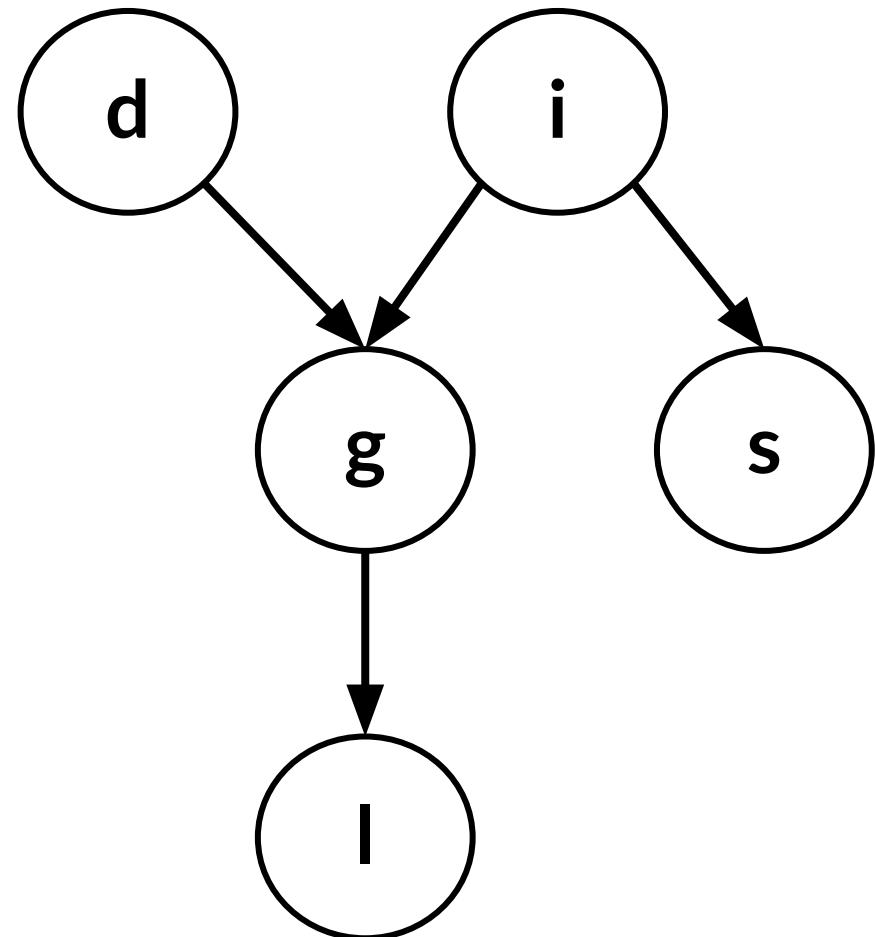
Can identify some conditional independencies by looking at graph properties.

Nodes in graph can correspond to random variables or random vectors.

Next: observed vs. latent variables; plate notation; famous Bayesian networks.

Bayesian Networks – Observed vs Latent

Back to the running example Bayesian network...



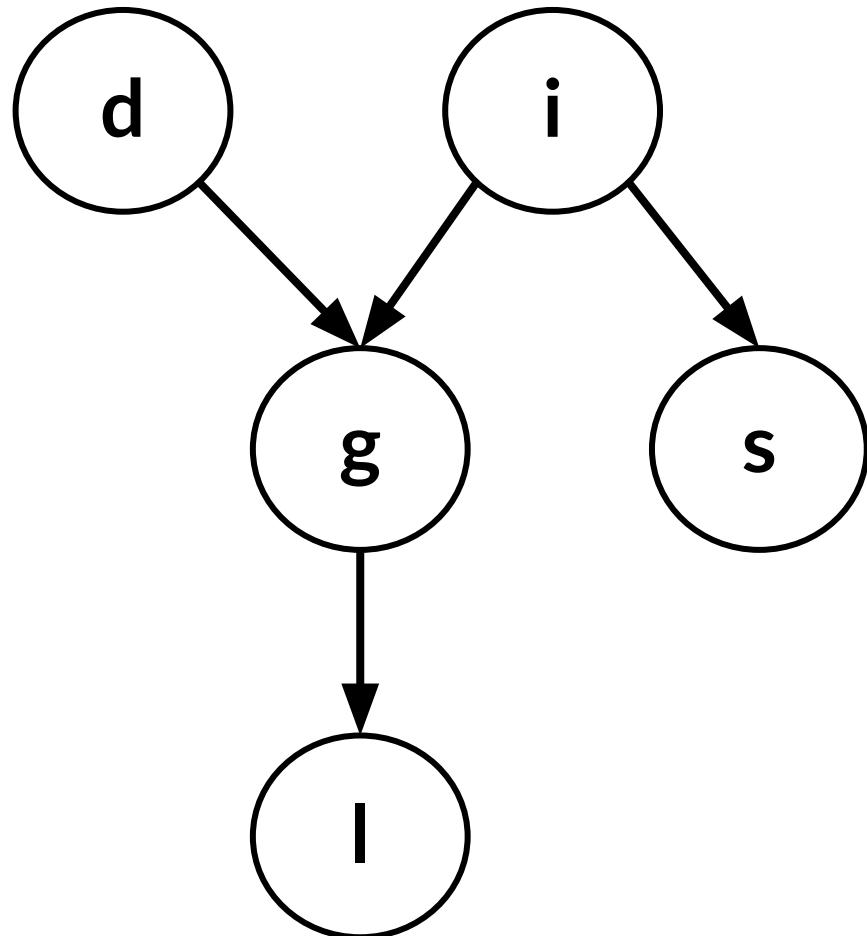
Bayesian Networks – Observed vs Latent

Observed Variables

- Variables that we observe (as samples from joint PDF)

Latent Variables

- Variables in the model that are unobserved.
- ⇒ but we may want to infer them given observed variables!



Bayesian Networks – Observed vs Latent

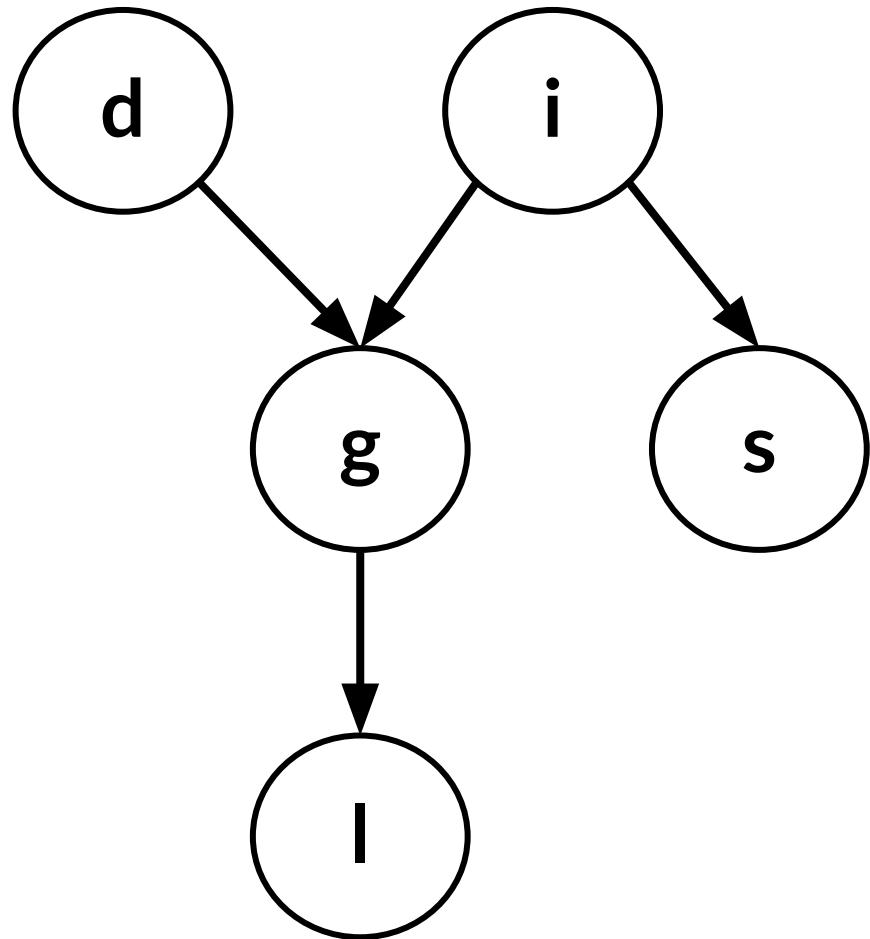
Observed Variables

- Variables that we observe (as samples from joint PDF)

Latent Variables

- Variables in the model that are unobserved.
- ⇒ but we may want to infer them given observed variables!

$$p(d | g, s, l)$$



Bayesian Networks – Observed vs Latent

Observed Variables

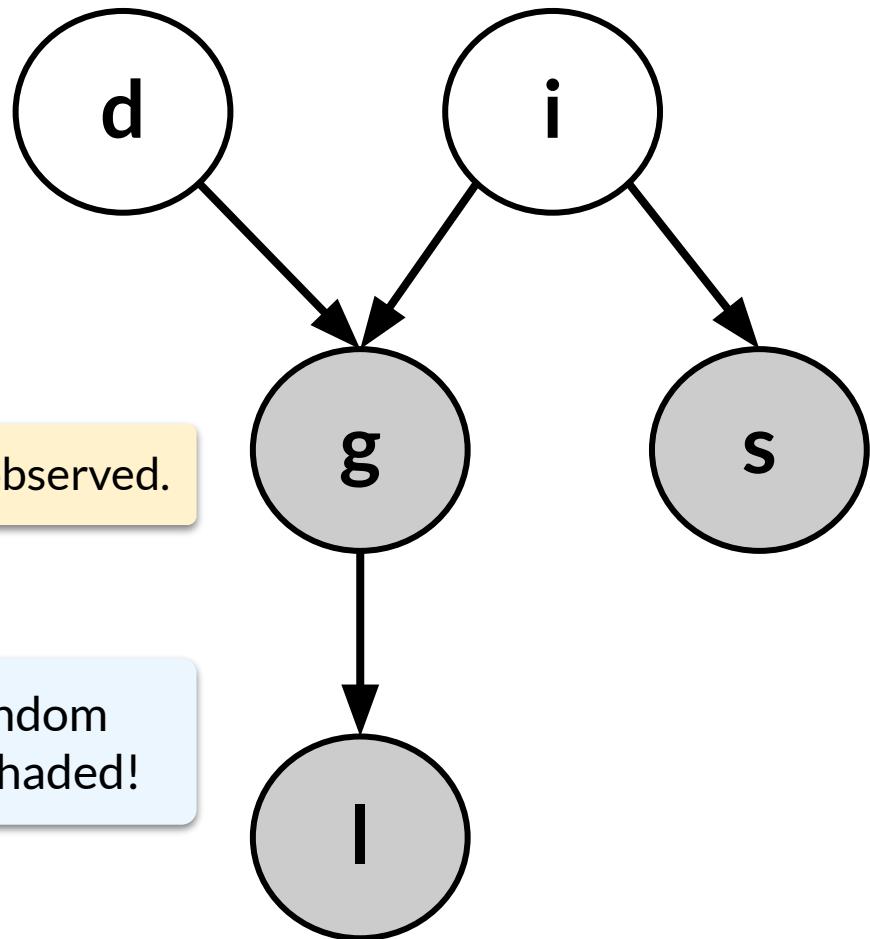
- Variables that we observe (as samples from joint PDF)

E.g., $\{d, i\}$ latent, $\{g, s, l\}$ observed.

Latent Variables

- Variables in the model that are unobserved.
- ⇒ but we may want to infer them given observed variables!

$$p(d | g, s, l)$$



Bayesian Networks – Observed vs Latent

Observed Variables

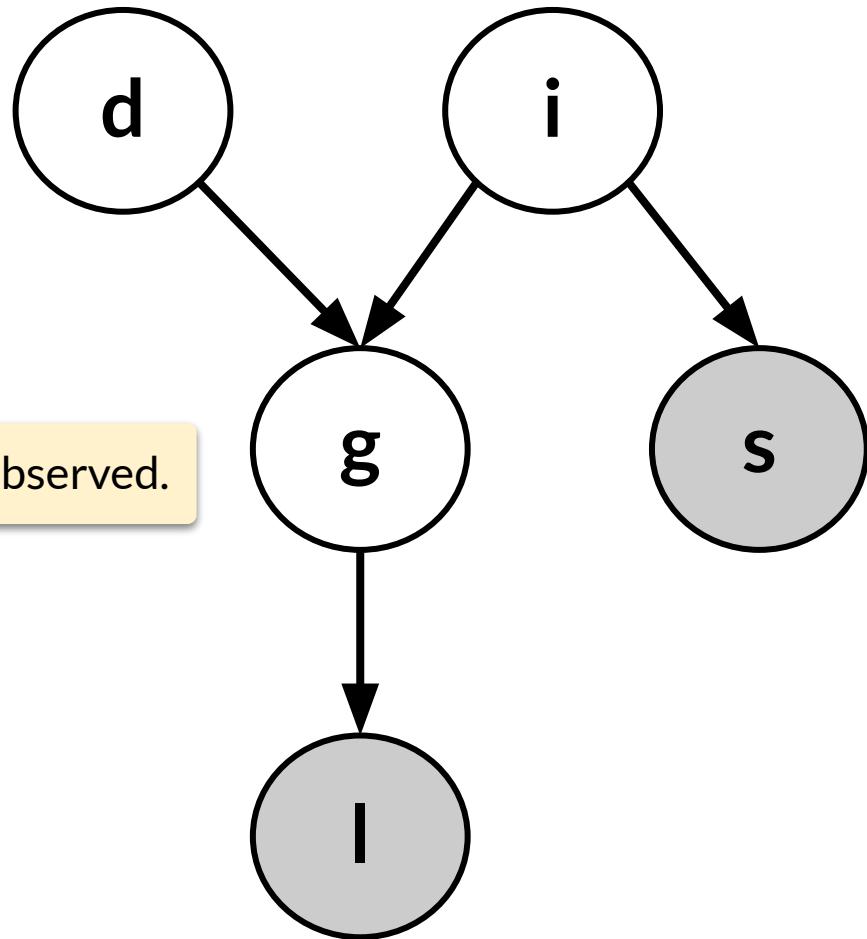
- Variables that we observe (as samples from joint PDF)

E.g., $\{d, i, g\}$ latent, $\{s, l\}$ observed.

Latent Variables

- Variables in the model that are unobserved.
- ⇒ but we may want to infer them given observed variables!

$$p(g \mid s, l)$$



Bayesian Networks – Observed vs Latent

Observed Variables

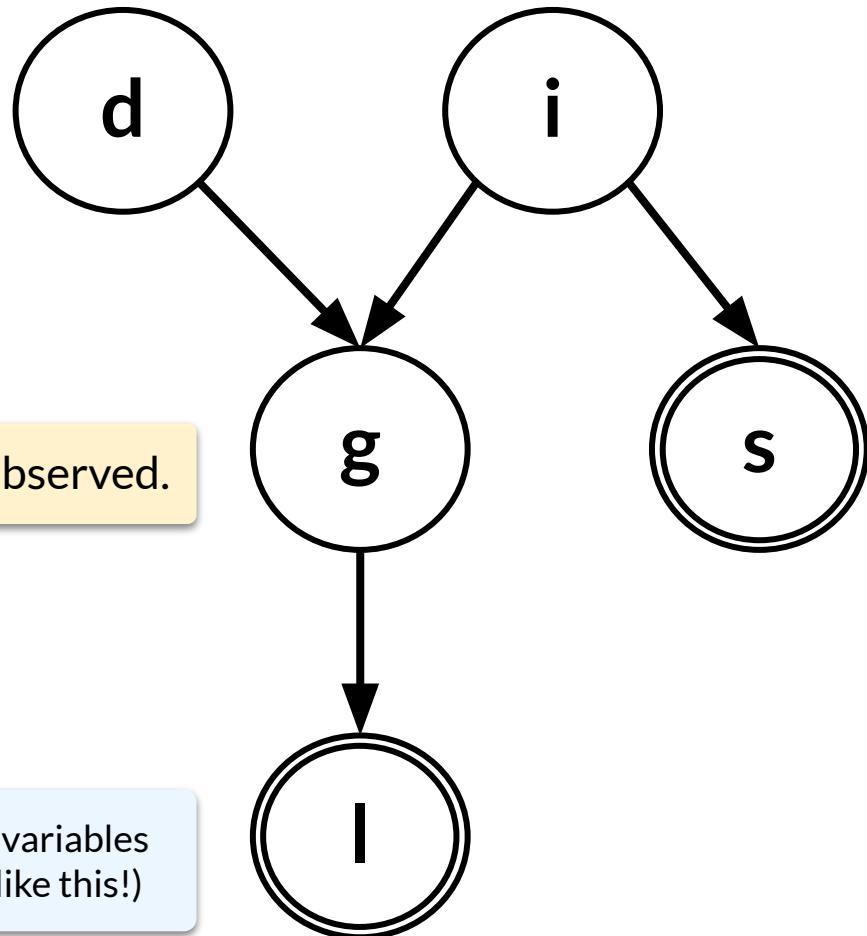
- Variables that we observe (as samples from joint PDF)

E.g., $\{d, i, g\}$ latent, $\{s, l\}$ observed.

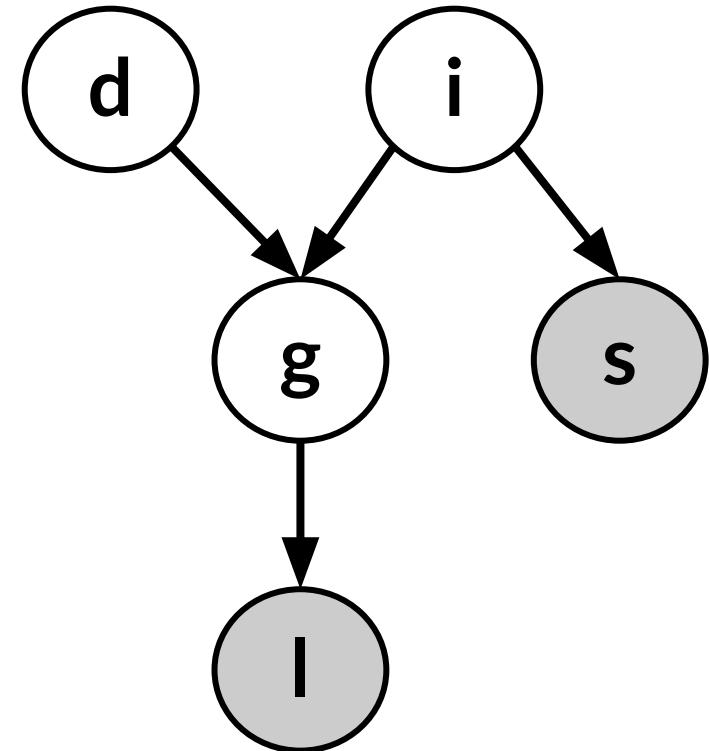
Latent Variables

- Variables in the model that are unobserved.
- ⇒ but we may want to infer them given observed variables!

$$p(g \mid s, l)$$

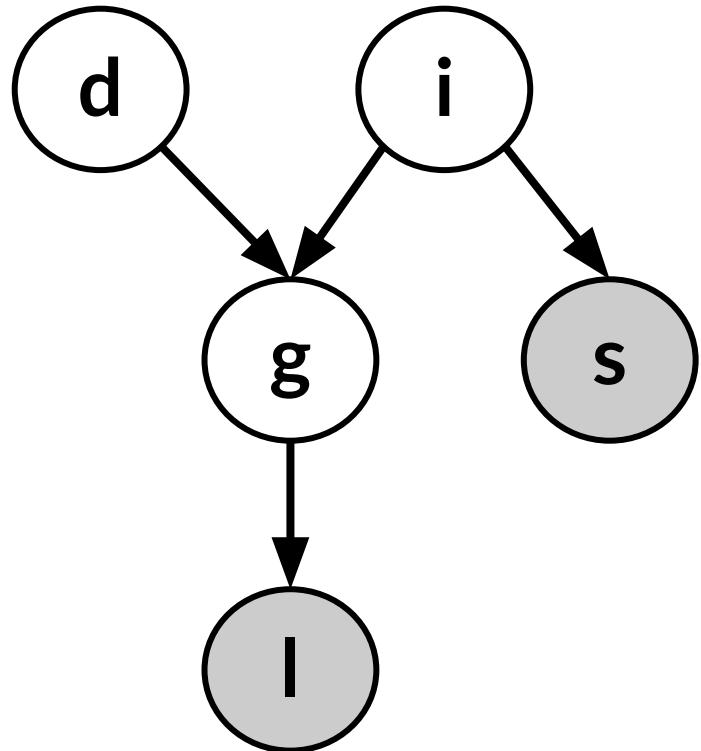


Bayesian Networks – Plate Notation



Bayesian Networks – Plate Notation

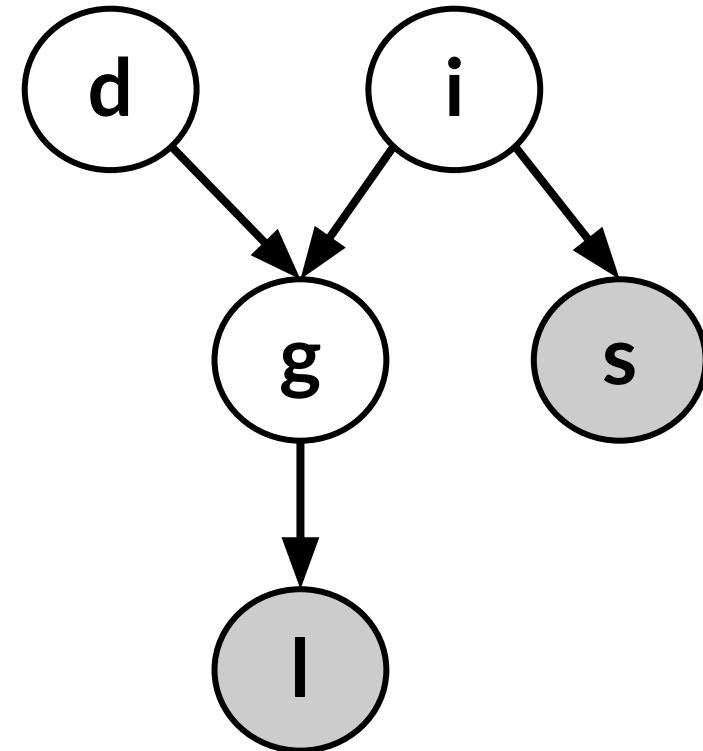
Plate notation is a “visual language” for describing more-complex Bayesian networks.



Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

⇒ Allows us to construct more-complex BNs with compressed notation!



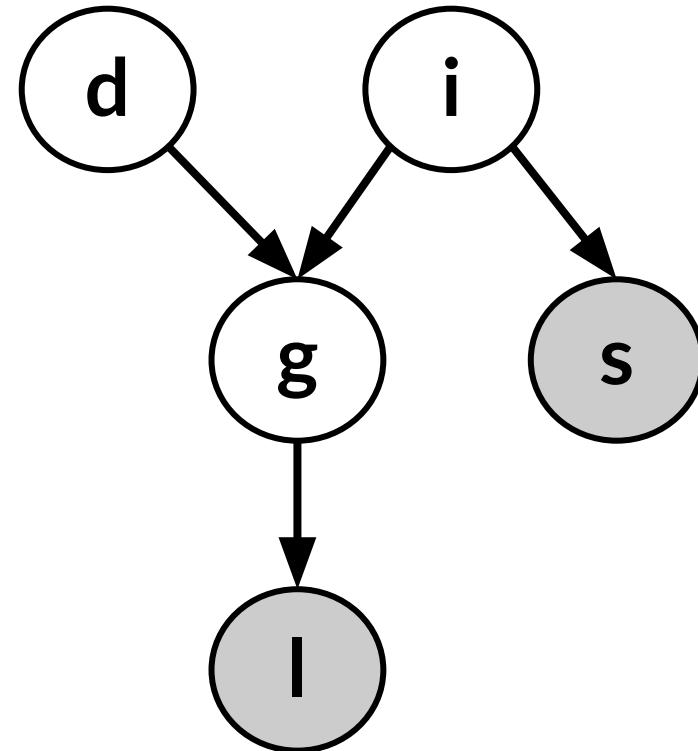
Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

⇒ Allows us to construct more-complex BNs with compressed notation!

E.g., what if we want to model 2 students?

(Each taking a separate class, an SAT test, and getting a letter-of-rec).



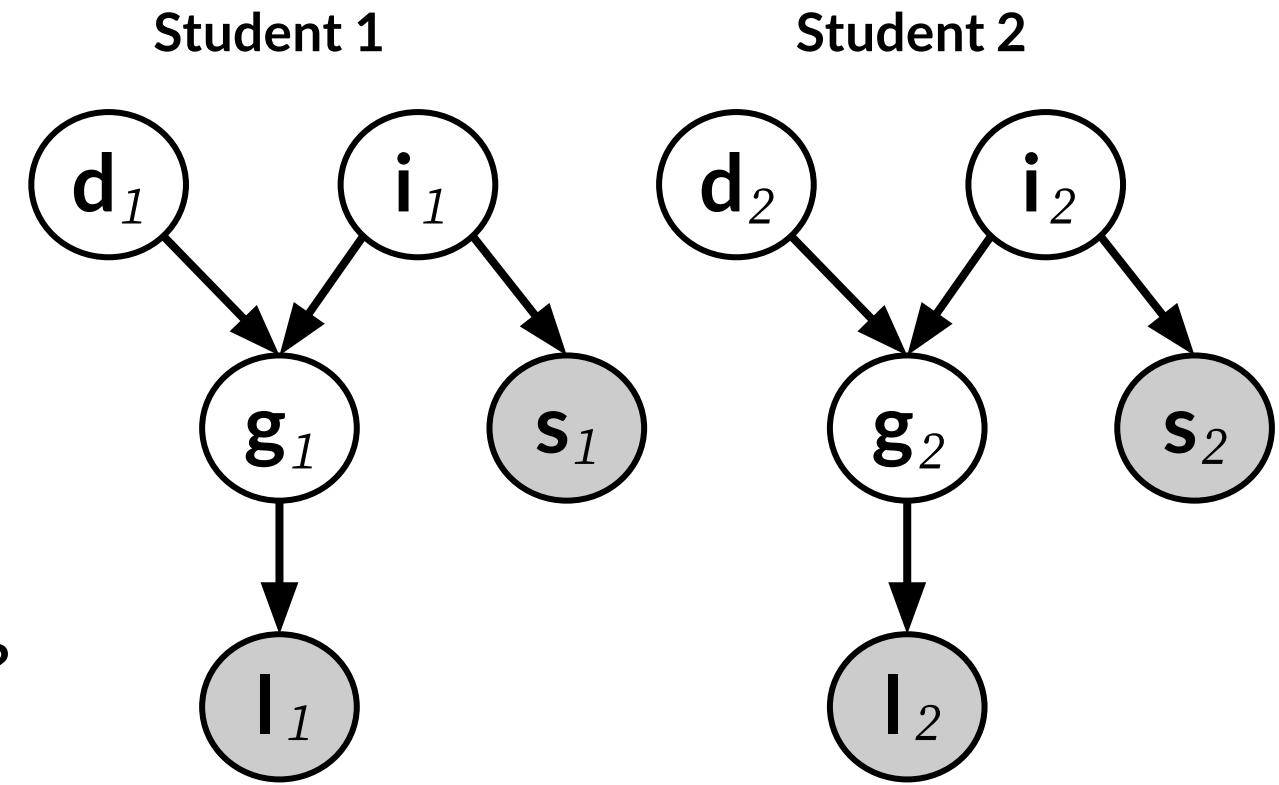
Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

⇒ Allows us to construct more-complex BNs with compressed notation!

E.g., what if we want to model 2 students?

(Each taking a separate class, an SAT test, and getting a letter-of-rec).



Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

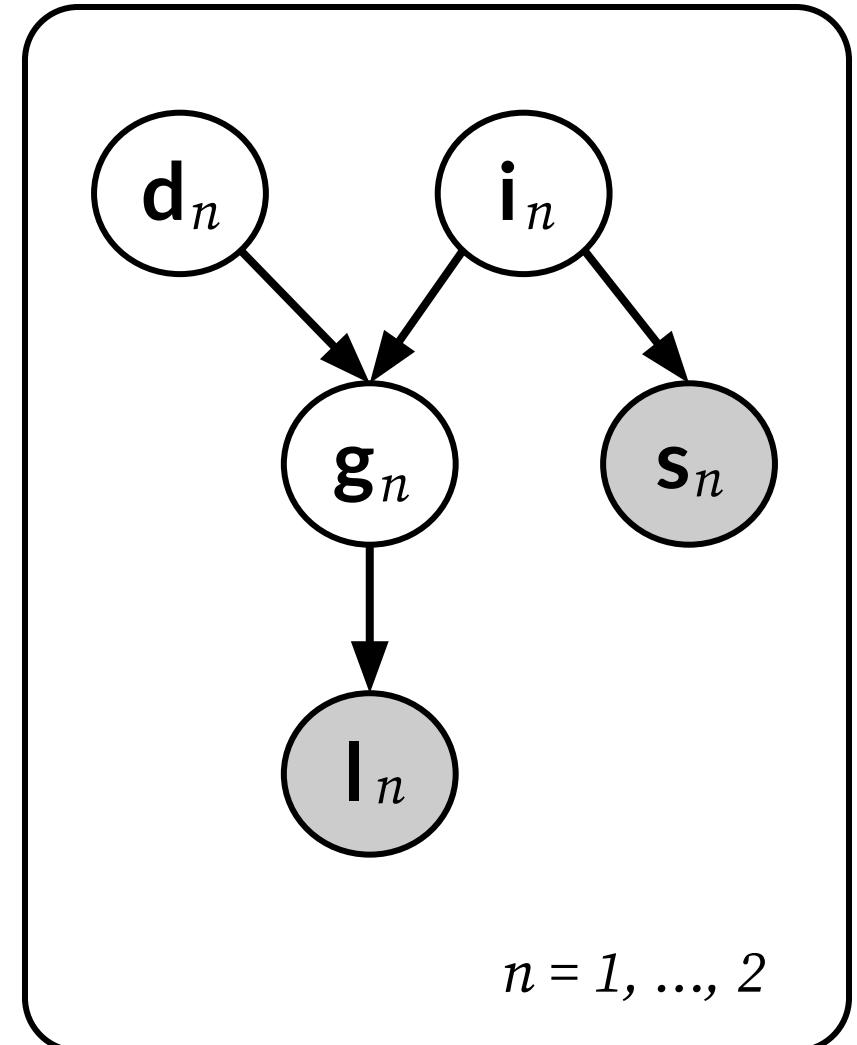
⇒ Plate with indices denotes sets of random variables.

Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

⇒ Plate with indices denotes sets of random variables.

Two (independent) sets of random variables.

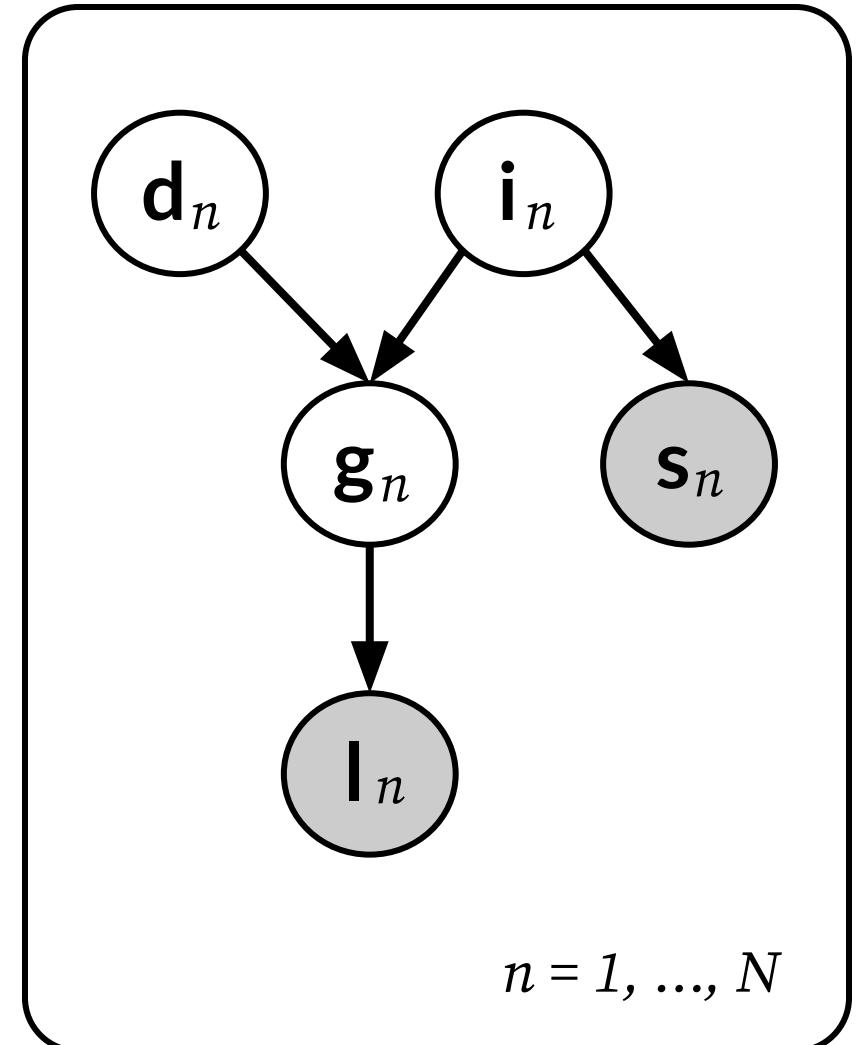


Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

⇒ Plate with indices denotes sets of random variables.

N (independent) sets of random variables.

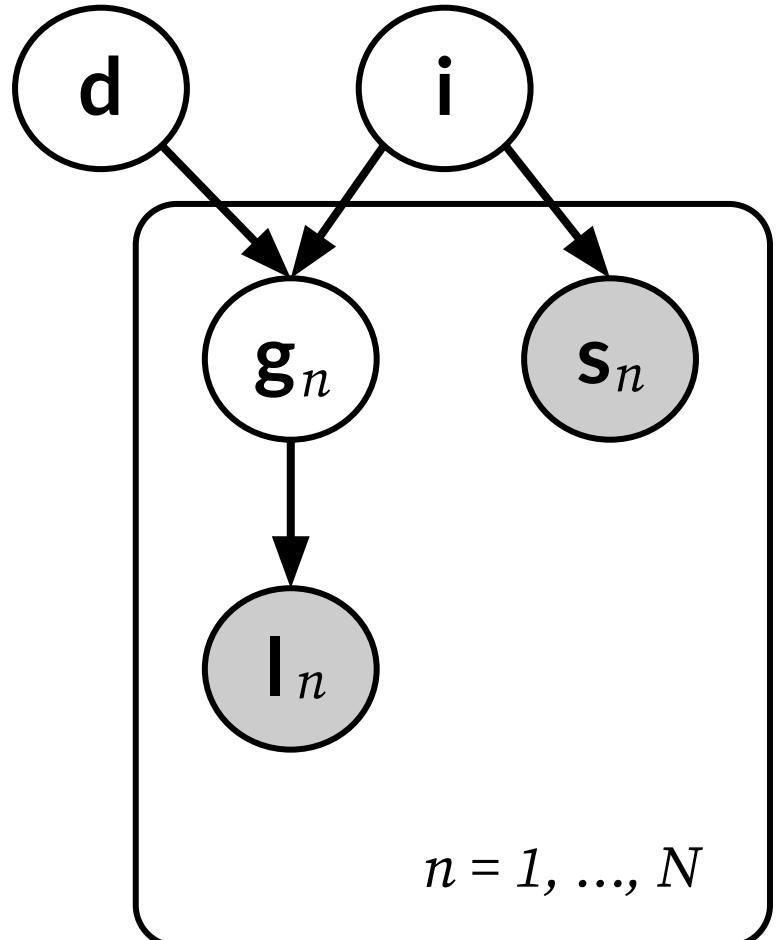


Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

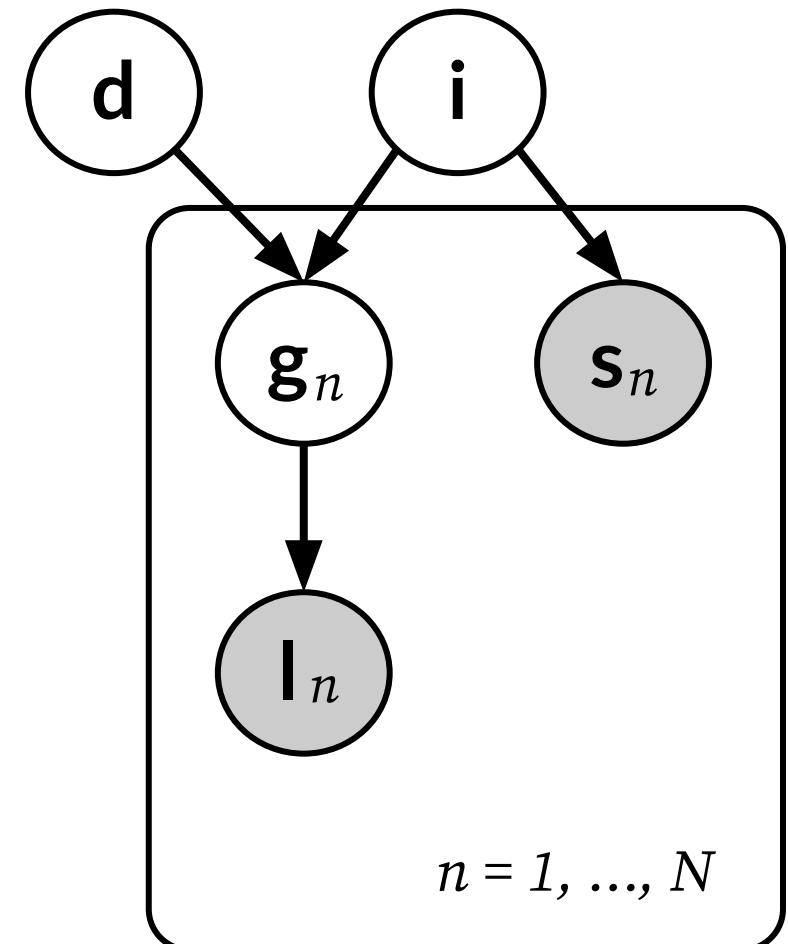
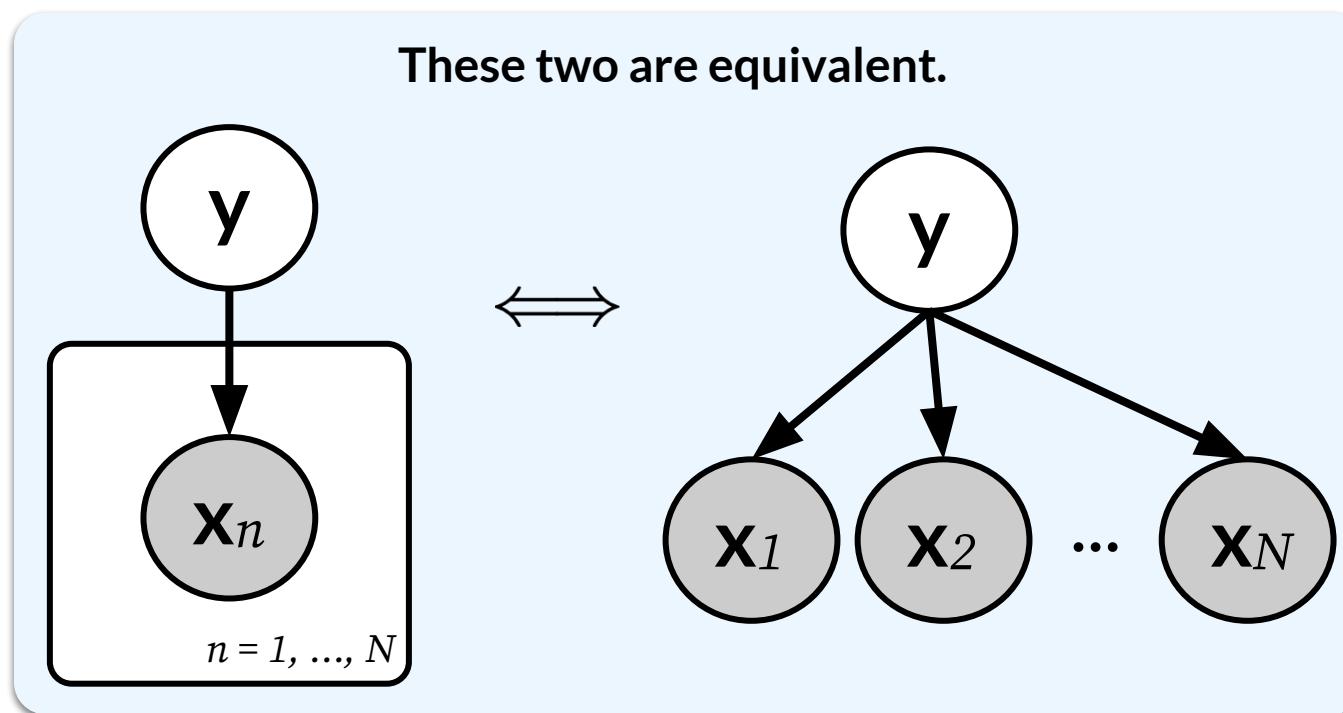
What does this plate mean?

(Where there are a subset of variables within the plate 🤔 ...)



Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

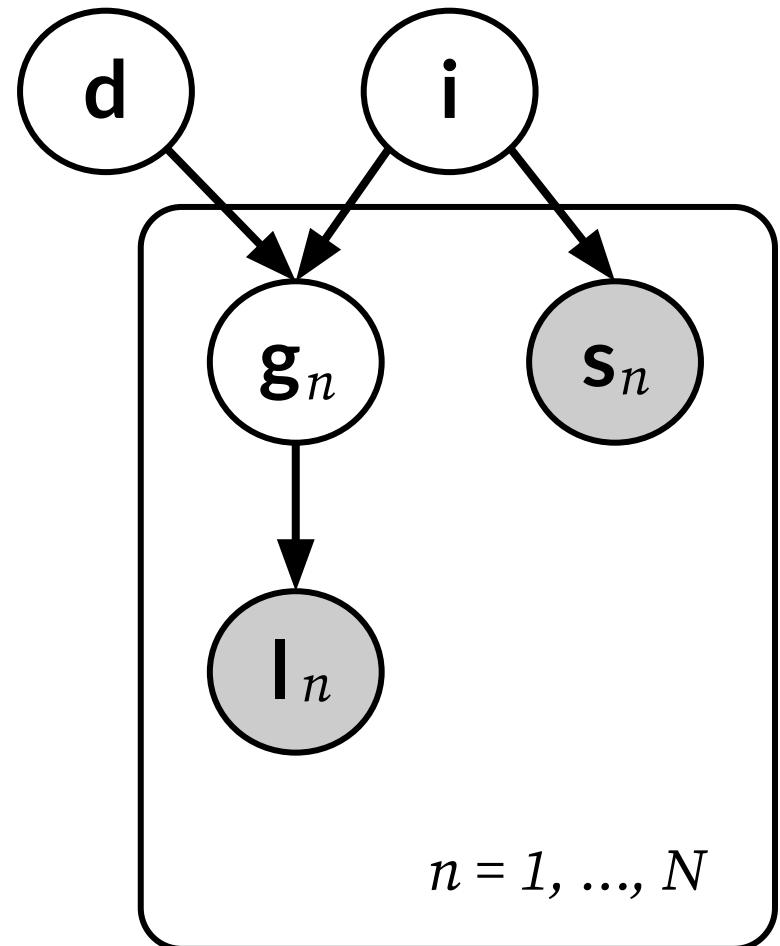


Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

This could be a model of (for example):

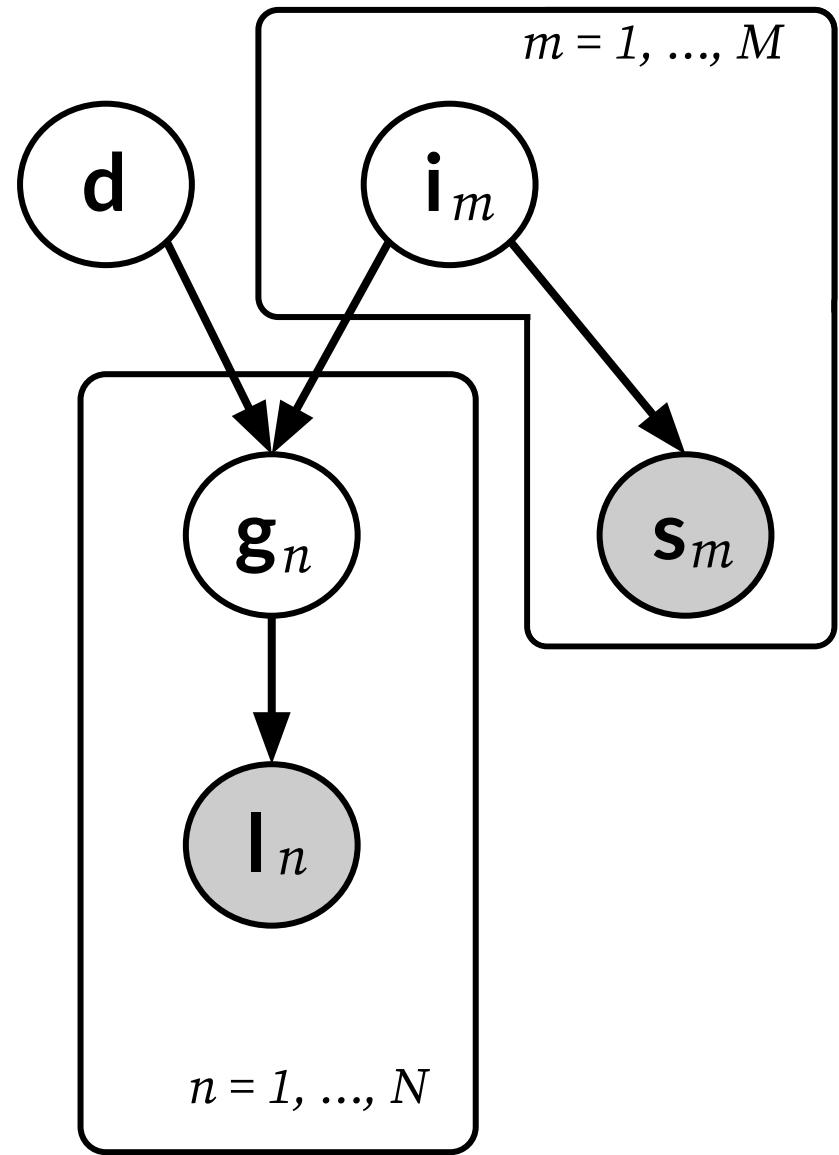
- A single student, with intelligence i , taking a single class with difficulty d .
- At N times during the semester, they:
 - Receive a grade g .
 - Get a letter of rec based on that grade.
 - Separately, take an SAT test.



Bayesian Networks – Plate Notation

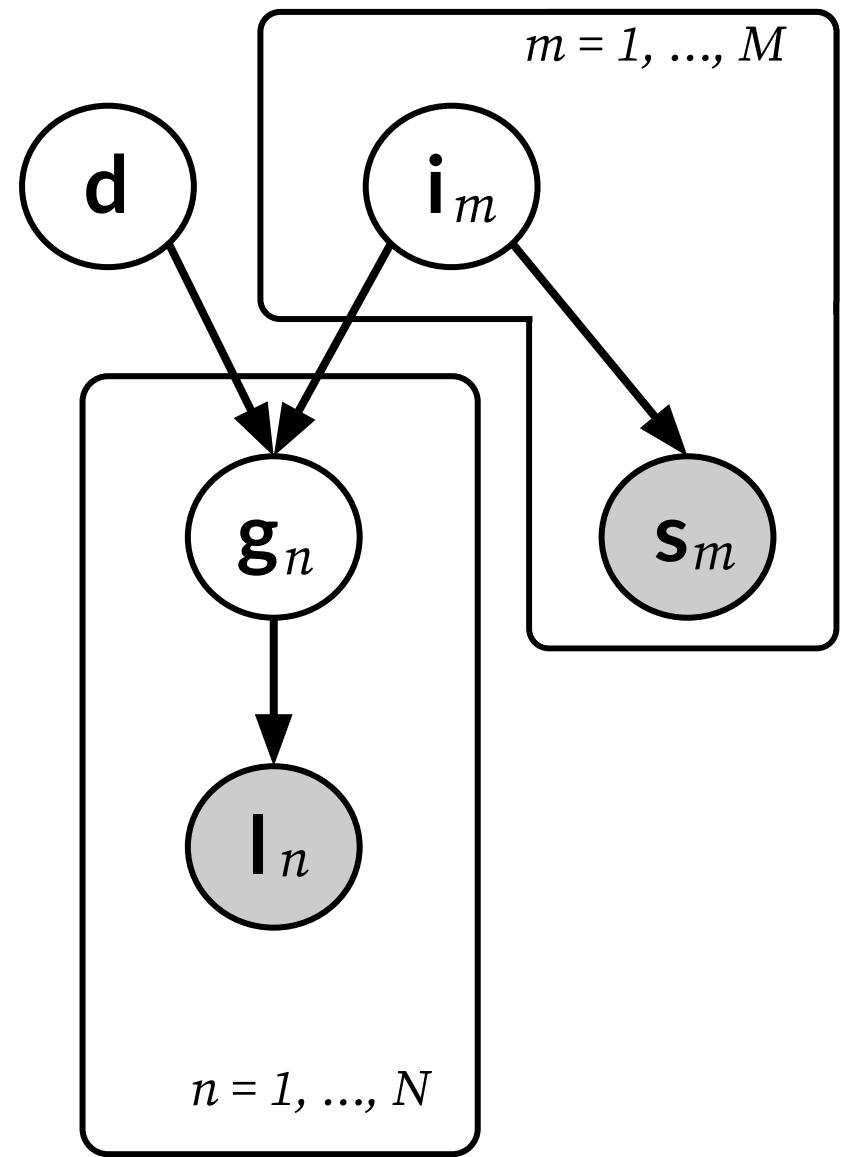
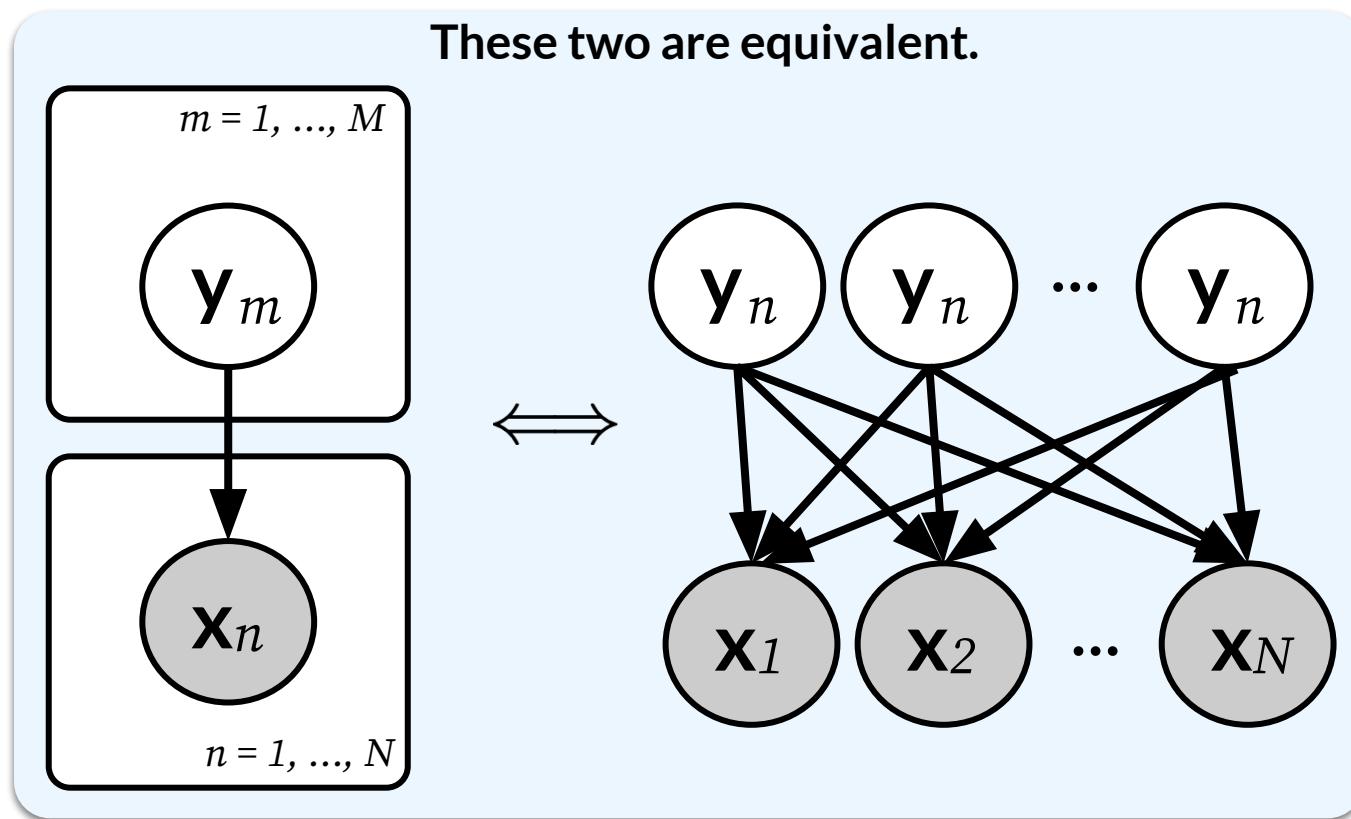
Plate notation is a “visual language” for describing more-complex Bayesian networks.

What about this structure?



Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

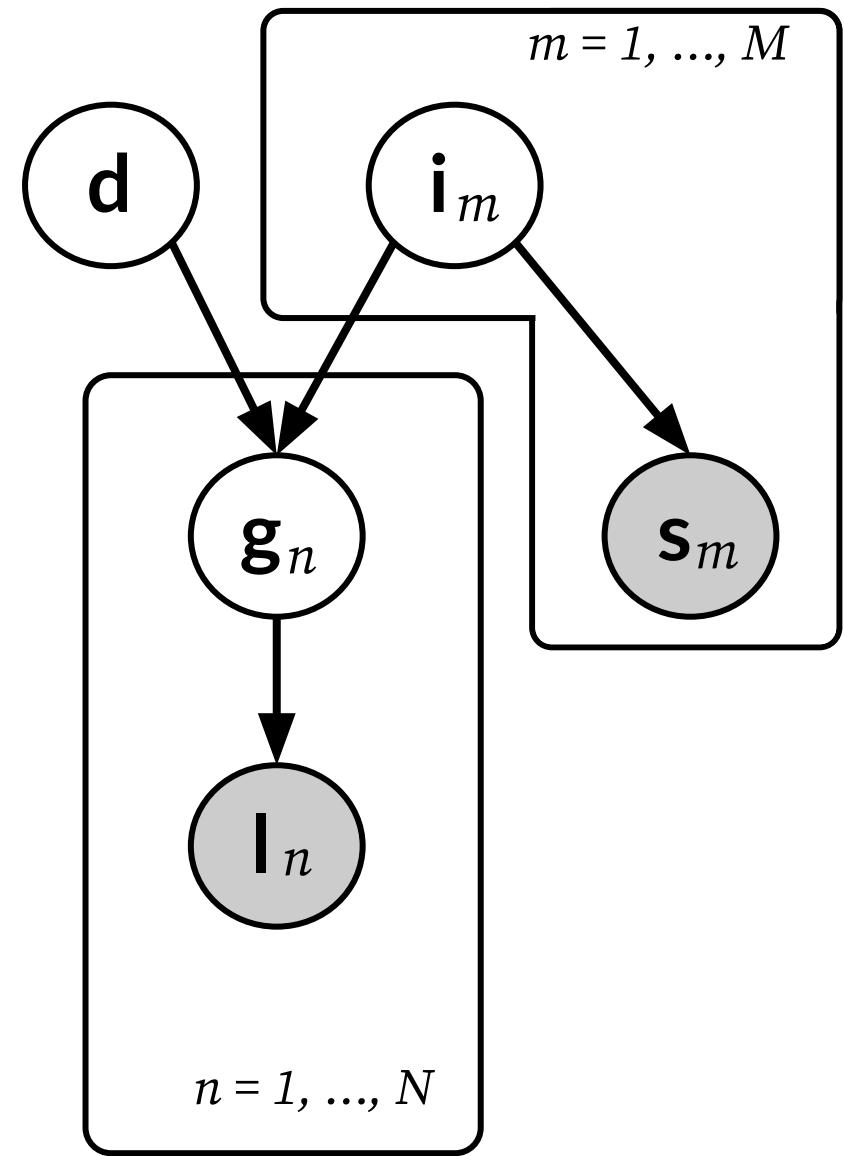


Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

This could be a model of (for example):

- M students, each with an intelligence i , in a group for a course project.
 - Each student in the group takes SAT test.
- At N times during the semester, the group:
 - Receives a grade g .
 - Gets a letter of rec based on that grade.

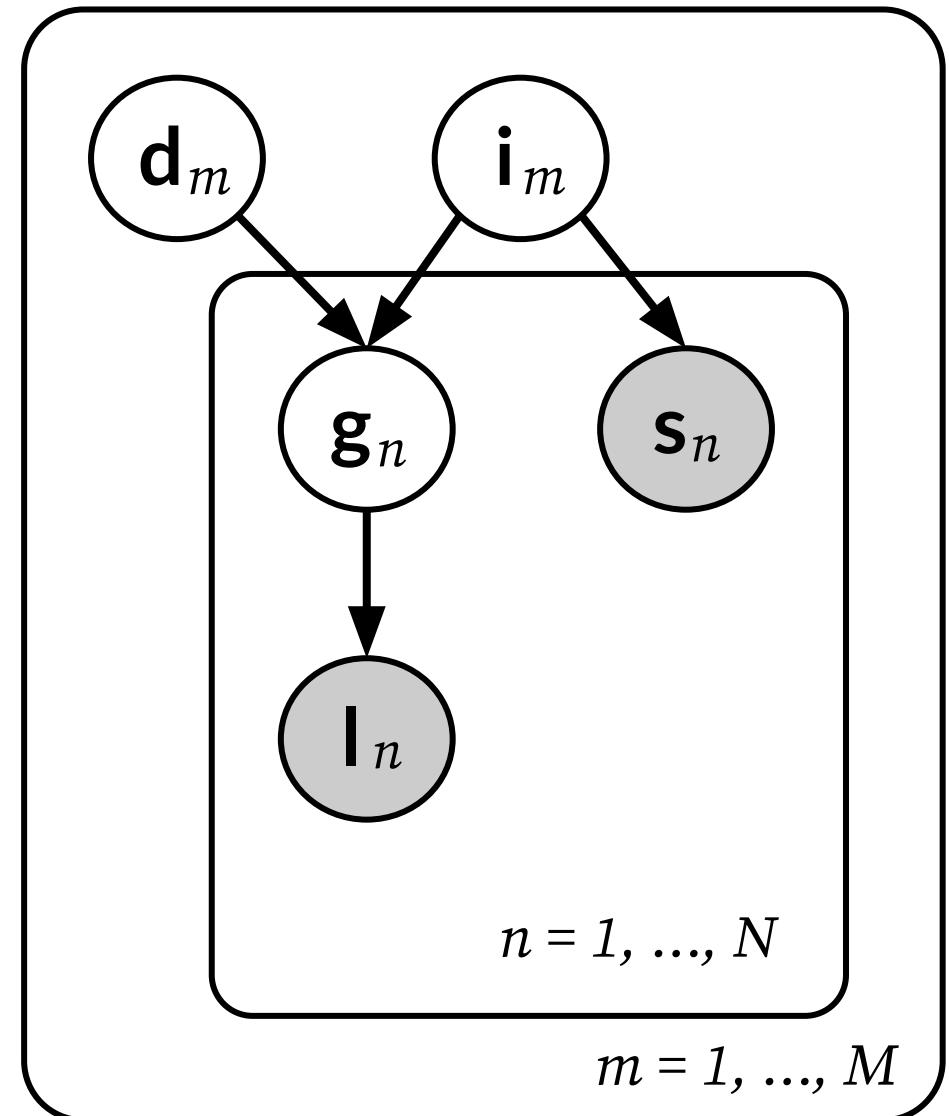


Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

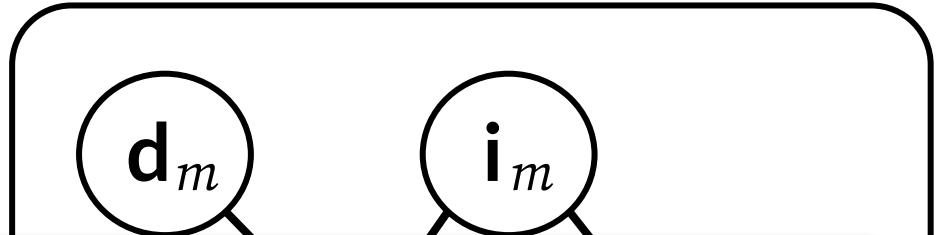
How about multiple (concentric) plates?

⇒ Yes!

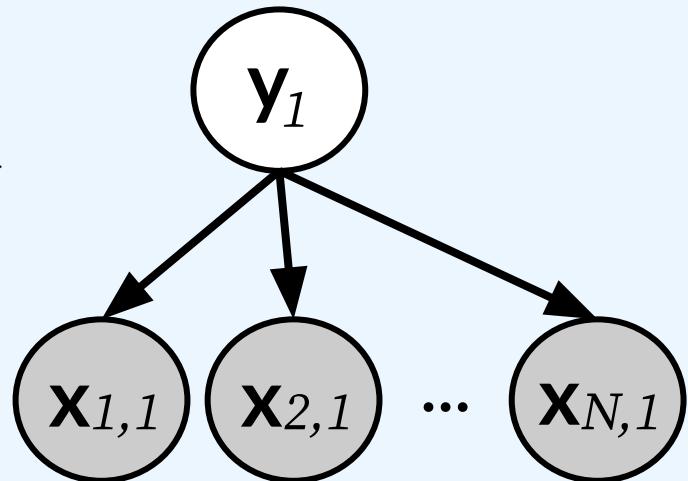
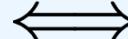
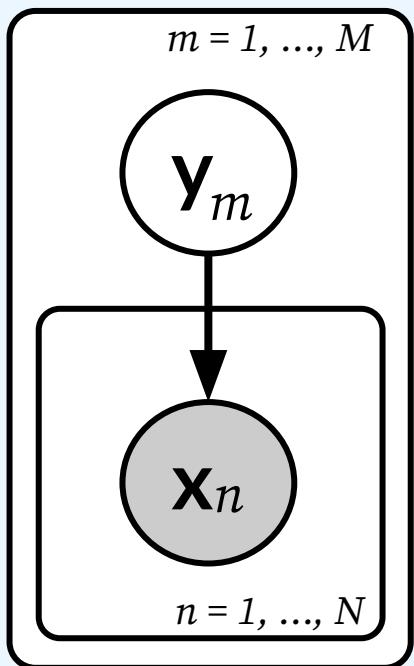


Bayesian Networks – Plate Notation

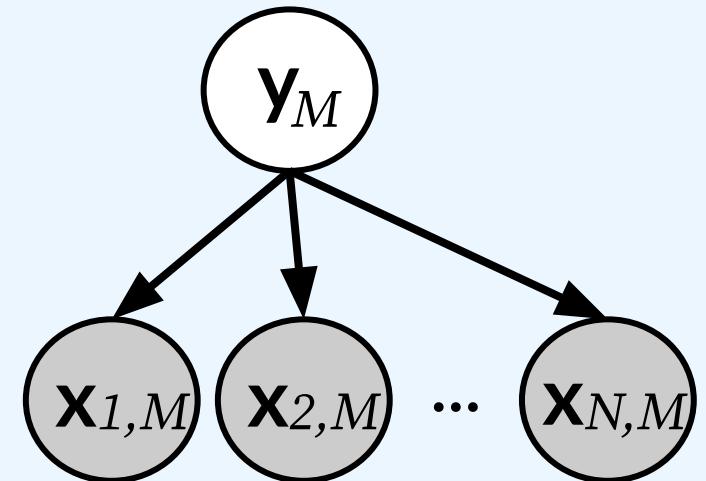
Plate notation is a “visual language” for describing more-complex Bayesian



These two are equivalent.



...



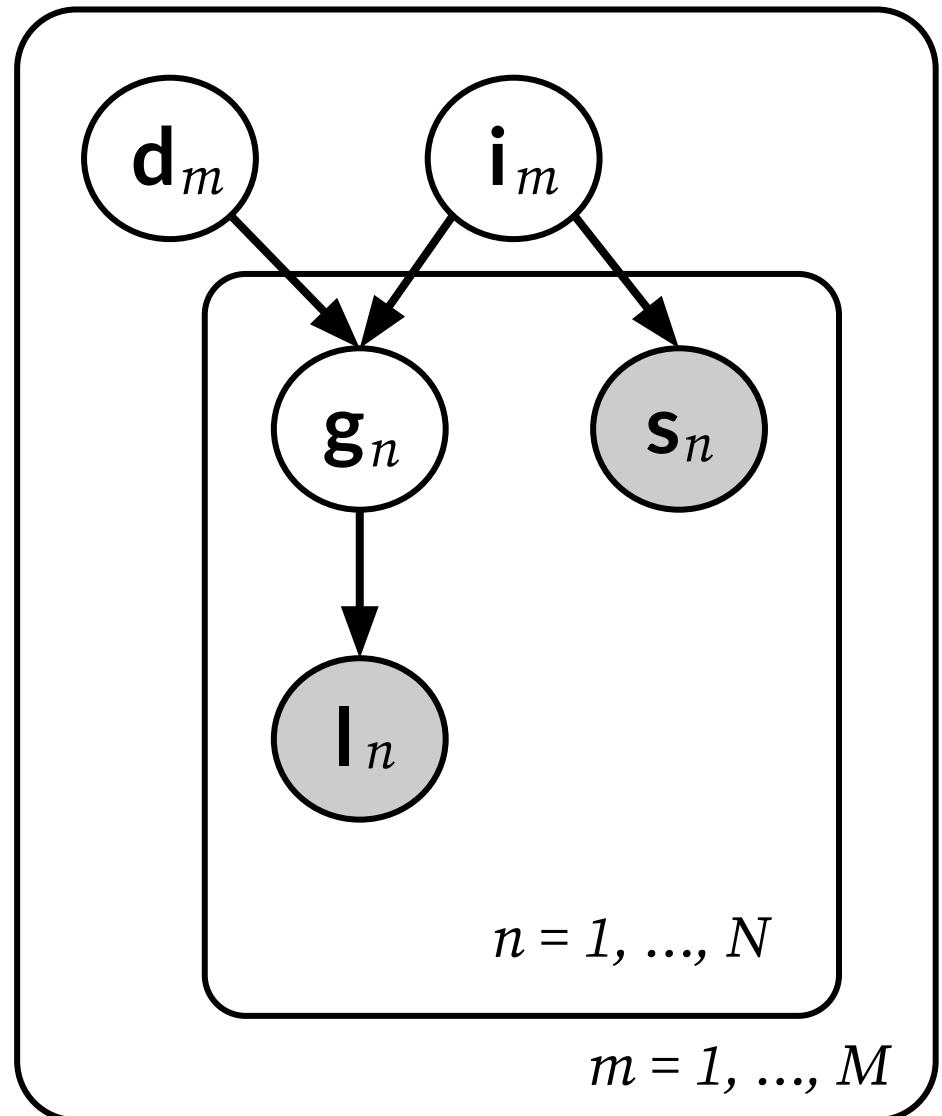
Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

Could be viewed as a model of...

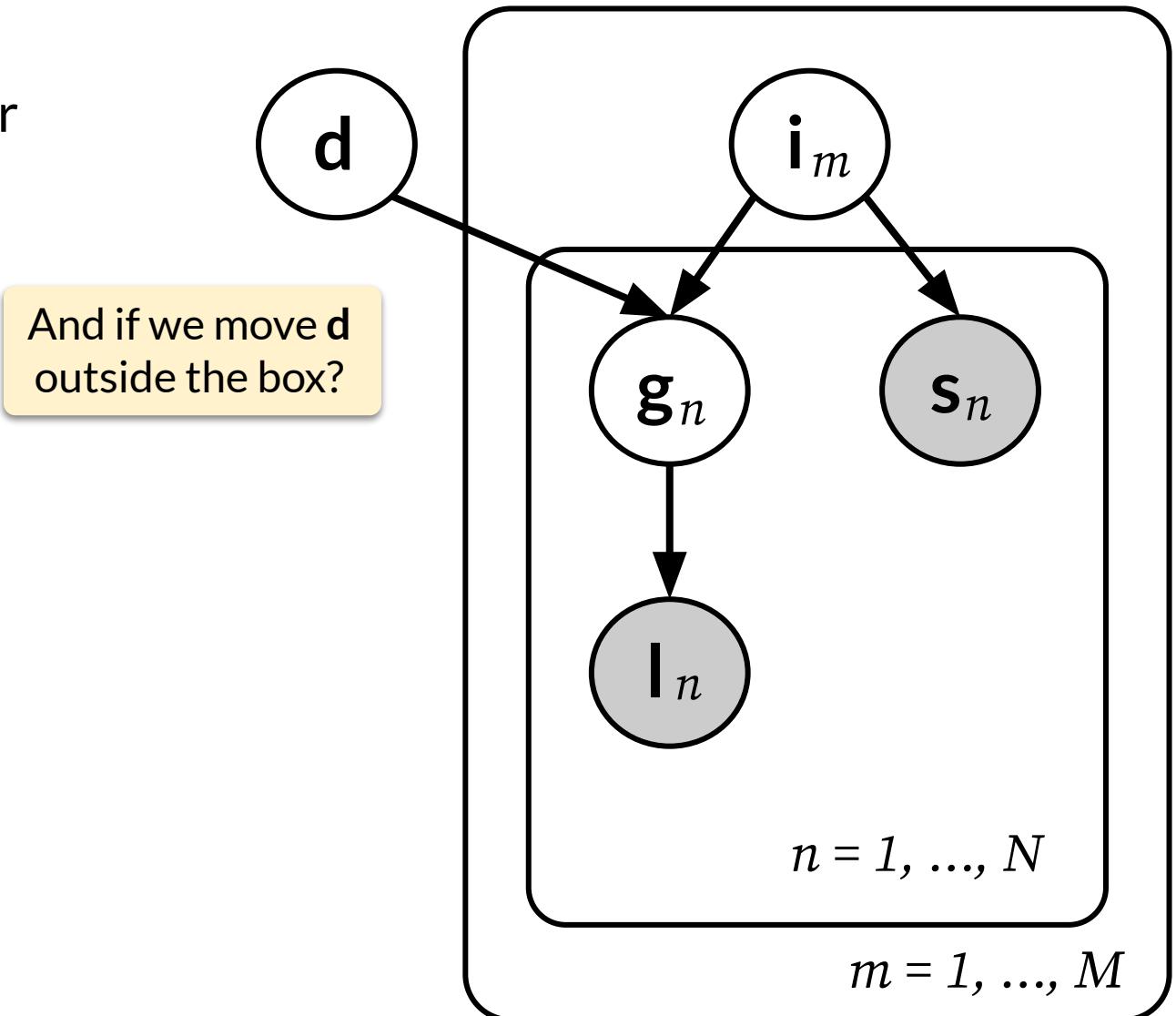
Repeat the following for M students:

- Each student, with intelligence i , is taking a (separate) class with difficulty d .
- At N times during the semester, they:
 - Receive a grade g .
 - Get a letter of rec based on that grade.
 - Separately, take an SAT test.



Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.



Bayesian Networks – Plate Notation

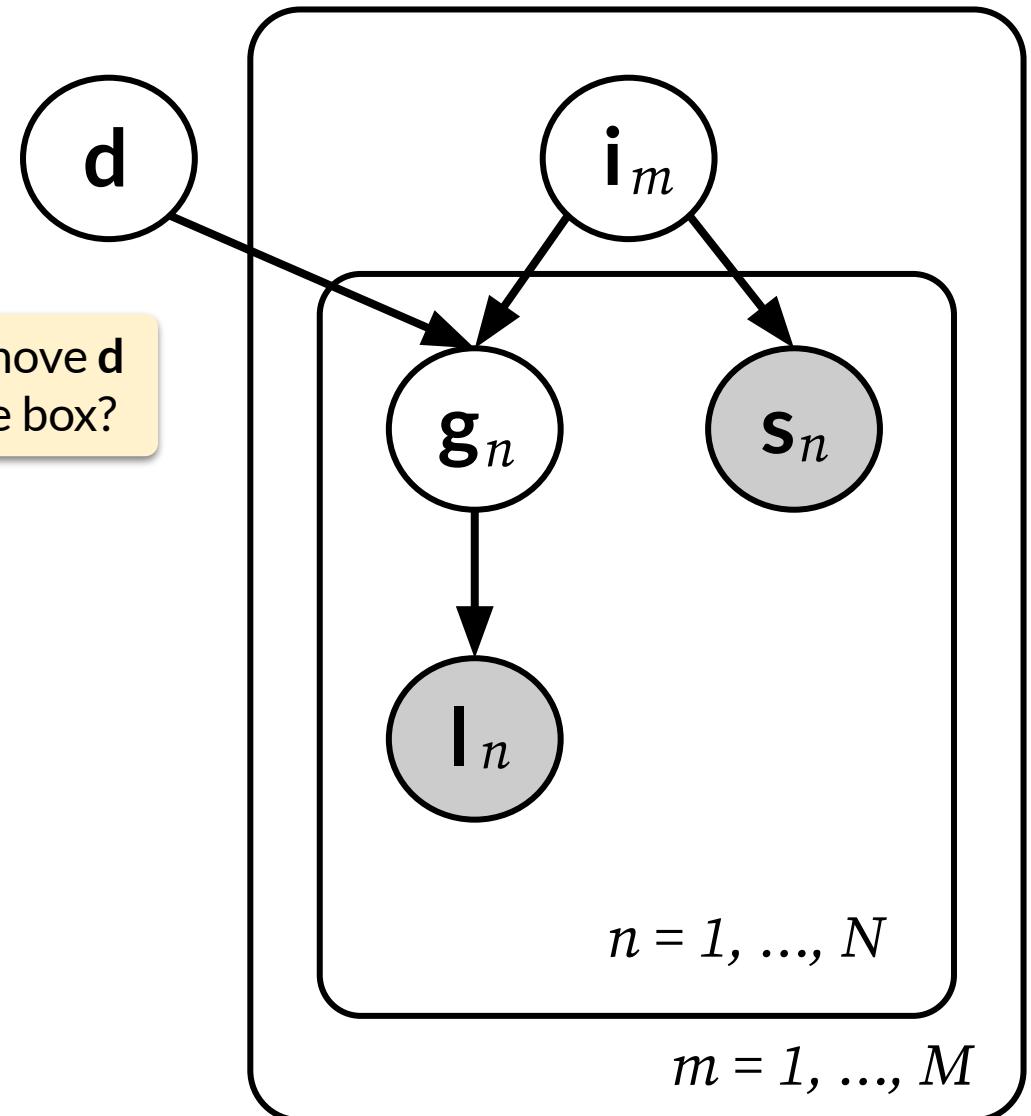
Plate notation is a “visual language” for describing more-complex Bayesian networks.

Could be viewed as a model of...

And if we move d outside the box?

Repeat the following for M students:

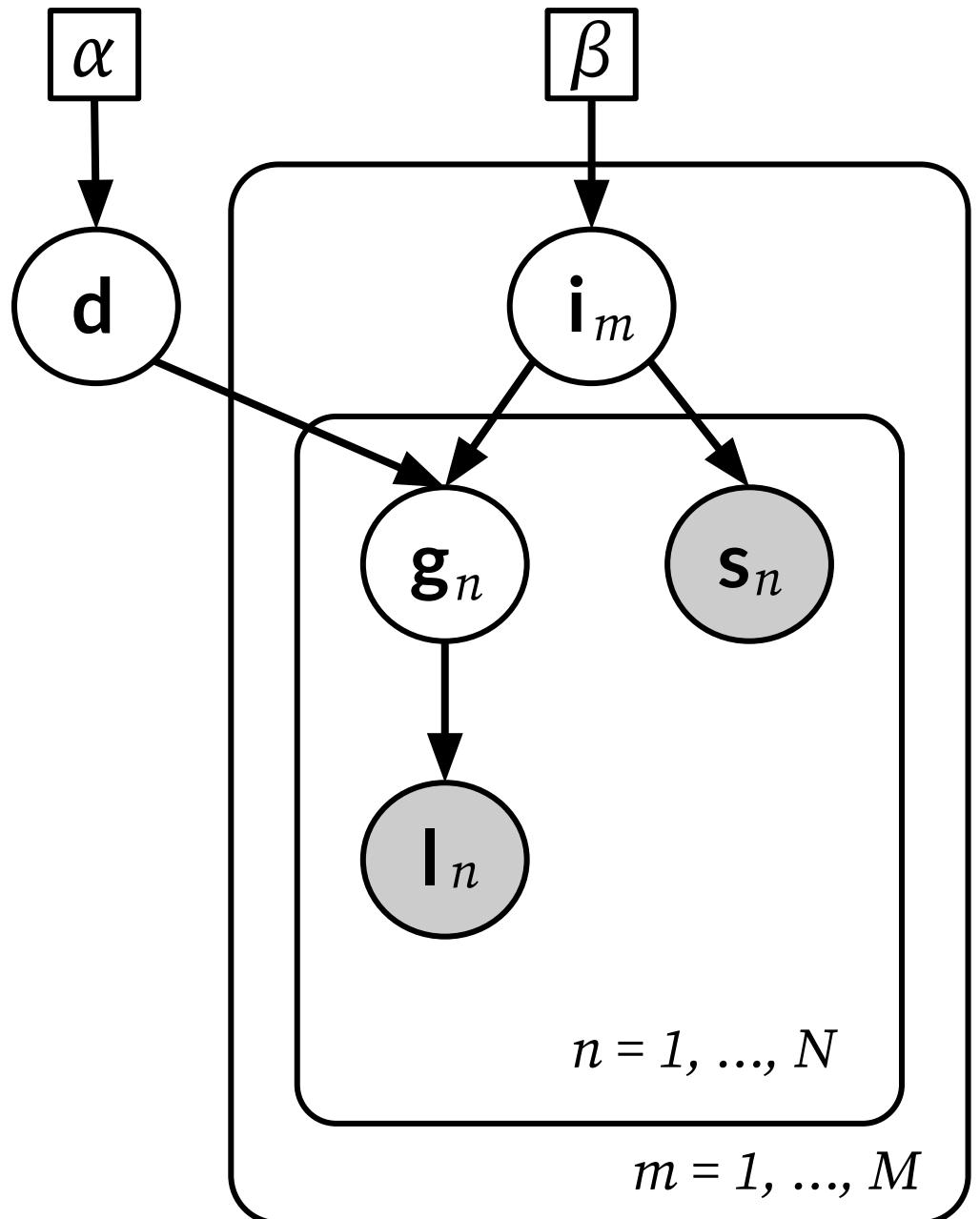
- A single student, with intelligence i , taking **the same** class with difficulty d .
- At N times during the semester, they:
 - Receive a grade g .
 - Get a letter of rec based on that grade.
 - Separately, take an SAT test.



Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

Sometimes, constant variables (i.e., **not random** variables) are included.



Bayesian Networks – Plate Notation

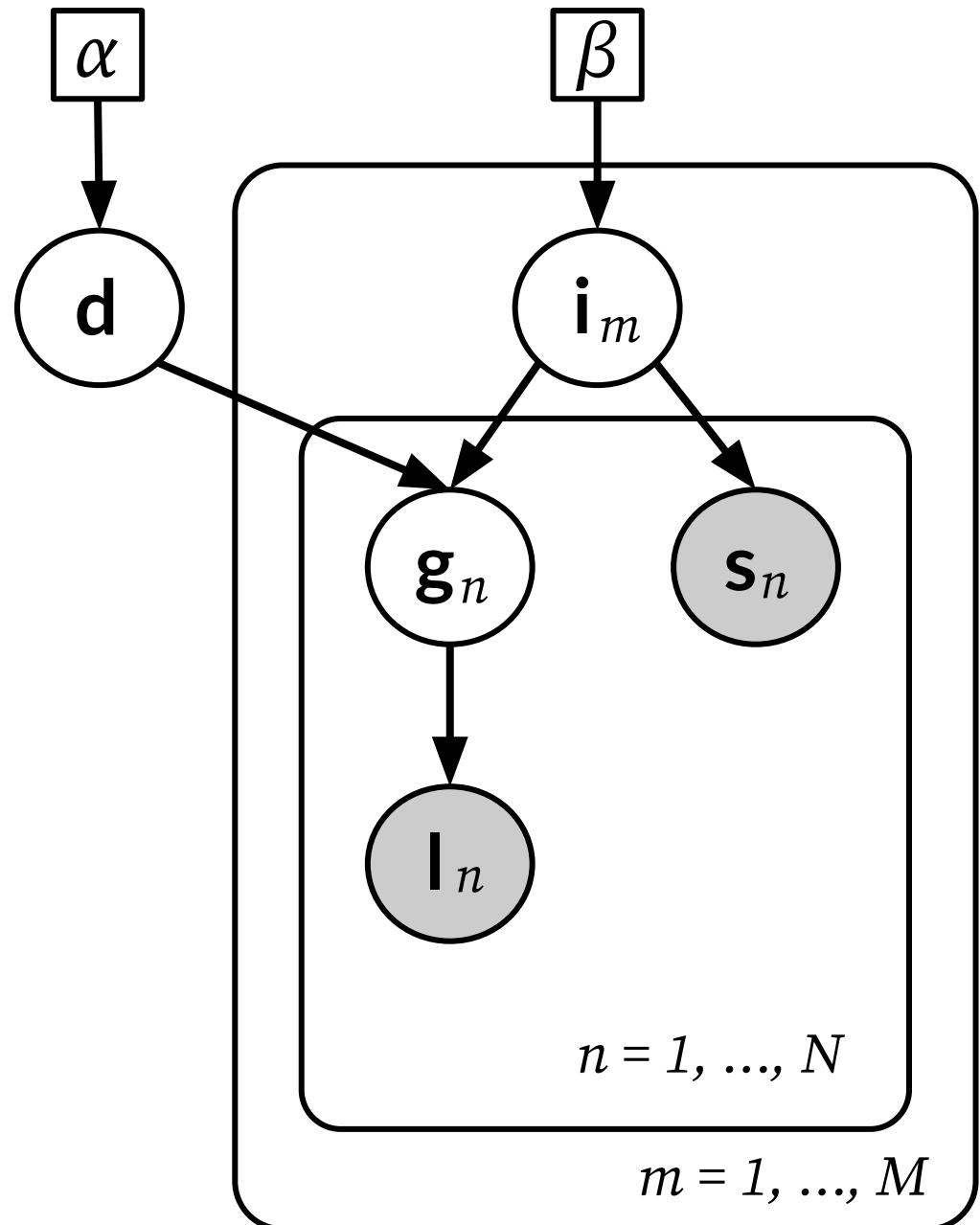
Plate notation is a “visual language” for describing more-complex Bayesian networks.

Sometimes, constant variables (*i.e., not random* variables) are included.

- For example, a parameter of a given distribution.
- (Often these are referred to as “hyperparameters”).
- Example:

$$d \sim \text{Gamma}(\alpha, 1)$$

$$i_m \sim \text{Bernoulli}(\beta), \text{ for } m = 1, \dots, M$$



Bayesian Networks – Plate Notation

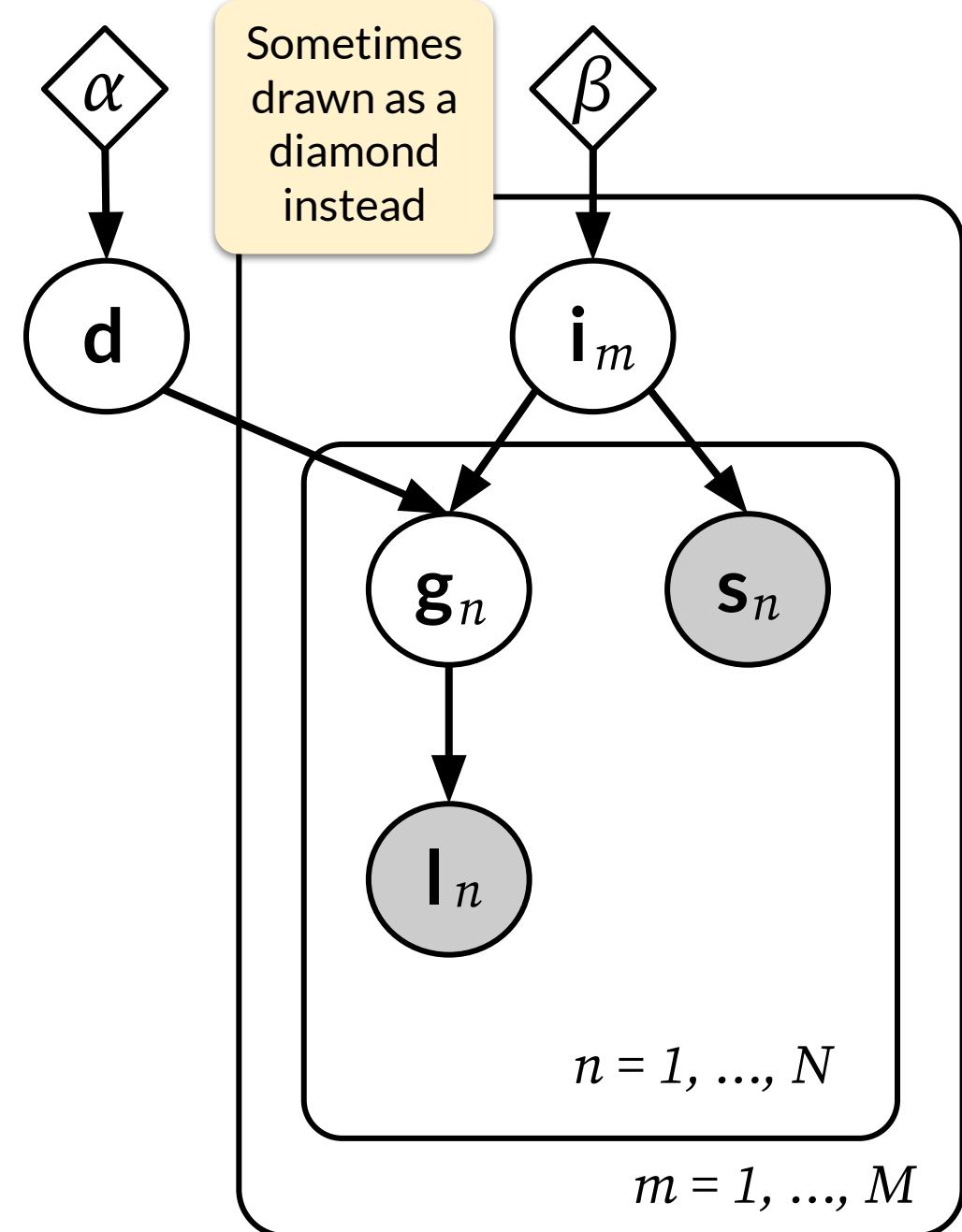
Plate notation is a “visual language” for describing more-complex Bayesian networks.

Sometimes, constant variables (i.e., **not random** variables) are included.

- For example, a parameter of a given distribution.
- (Often these are referred to as “hyperparameters”).
- Example:

$$d \sim \text{Gamma}(\alpha, 1)$$

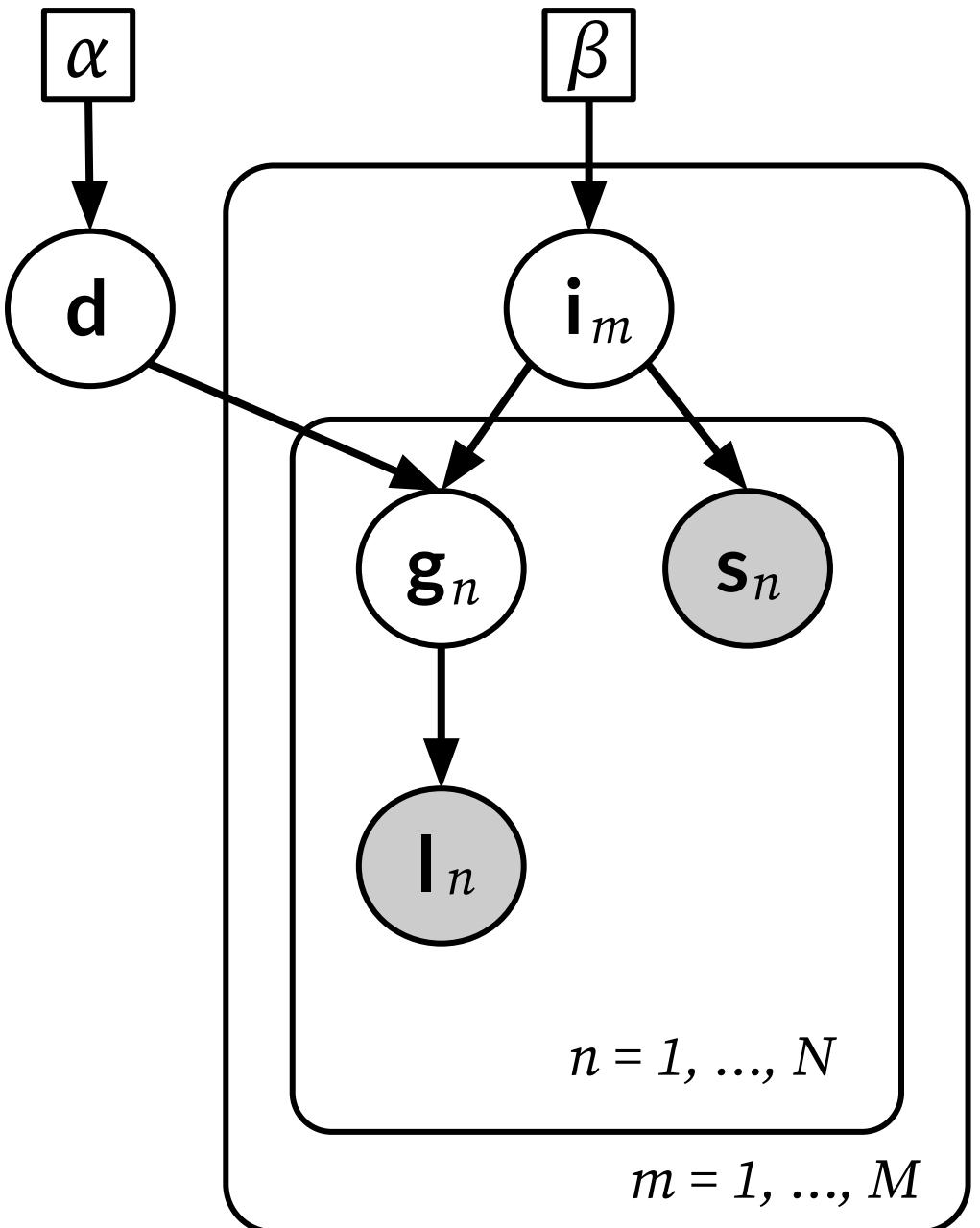
$$i_m \sim \text{Bernoulli}(\beta), \text{ for } m = 1, \dots, M$$



Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

To summarize, main rules of plate notation:

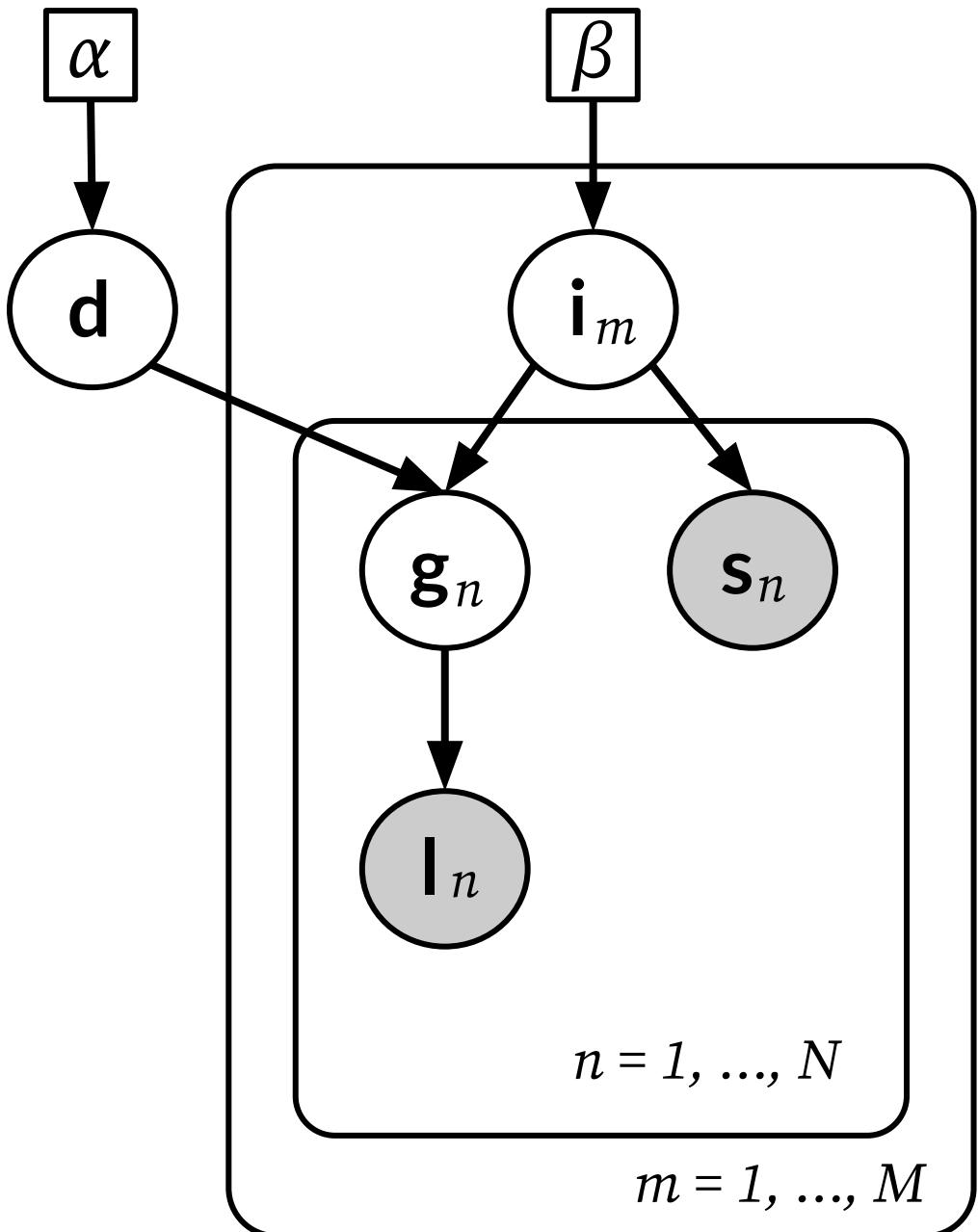


Bayesian Networks – Plate Notation

Plate notation is a “visual language” for describing more-complex Bayesian networks.

To summarize, main rules of plate notation:

- Each random variable (or vector) has a node.
- Shaded=observed variable. Unshaded=latent variable.
- Constants (eg, hyperparameters) as squares/diamonds.
- Plates with indices denote repeated variables.
- Variables within a plate do not have directed edges between them.



Bayesian Networks – Summary of Benefits

To summarize.... *What are the benefits of representing a joint PDF in this way (as a PGM)?*

Bayesian Networks – Summary of Benefits

To summarize.... *What are the benefits of representing a joint PDF in this way (as a PGM)?*

Can define & identify conditional independencies (reductions in complexity) via the graph structure.

Bayesian Networks – Summary of Benefits

To summarize.... *What are the benefits of representing a joint PDF in this way (as a PGM)?*

Can define & identify conditional independencies (reductions in complexity) via the graph structure.

Leads to algorithms for inference (computing conditionals/marginals) **based on the graph structure.**

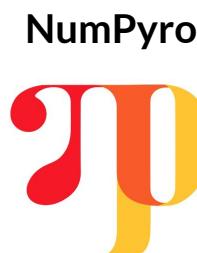
Bayesian Networks – Summary of Benefits

To summarize.... *What are the benefits of representing a joint PDF in this way (as a PGM)?*

Can define & identify conditional independencies (reductions in complexity) via the graph structure.

Leads to algorithms for inference (computing conditionals/marginals) **based on the graph structure**.

Sometimes this graph/plate data structure is used to specify a probabilistic model for deployment of automatic inference in probabilistic programming languages.



NumPyro

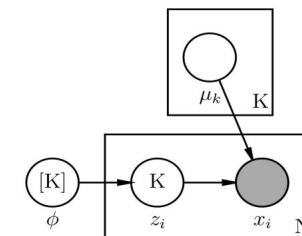


Plate notation of Gaussian Mixture Model

```
def gmm(data, K):
    phi = sample("phi", Dirichlet(np.ones(K)))
    with plate("K", K, dim=-1):
        mu = sample("mu", Normal(np.arange(K), 1))
    with plate("N", len(data), dim=-1):
        z = sample("z", Categorical(phi))
        sample("obs", Normal(mu[z], 1), obs=data)
```

Preview – Next Class

In next class, we will go into the details of inference algorithms (for computing the conditionals/marginals) based on the graph structure!

But first in this class, want to show some example (famous) Bayesian networks, and give some examples of where they've been used...

Famous Bayesian Networks!

Going through examples of famous Bayesian networks (and what they model).

- Markov models (Markov chains).
- Hidden Markov models (HMMs).
- Gaussian mixture models (GMMs).
- Latent Dirichlet allocation (LDA).
- Variational Autoencoder (VAE).

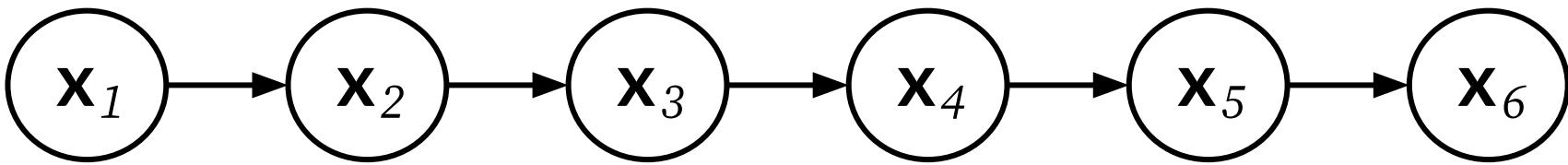
Famous Bayesian Networks – Markov Models

Famous Bayesian Networks – Markov Models

Markov Models ... also known as **Markov chains**.

Famous Bayesian Networks – Markov Models

Markov Models ... also known as **Markov chains**.



Generative process:

$$x_1 \sim p(\cdot)$$

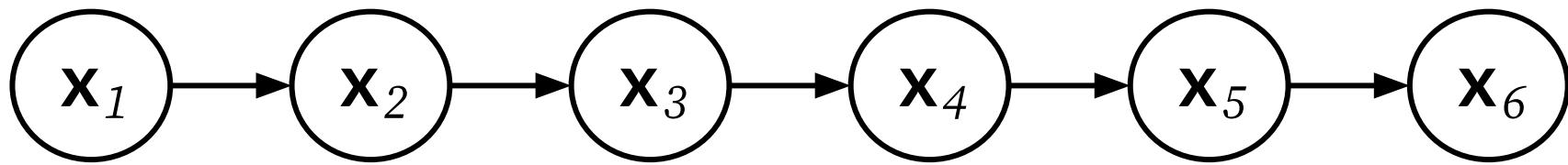
$$x_2 \sim p(x_1)$$

...

$$x_n \sim p(x_{n-1})$$

Famous Bayesian Networks – Markov Models

Markov Models ... also known as **Markov chains**.



Generative process:

$$x_1 \sim \mathcal{N}(0, 1)$$

$$x_2 \sim \mathcal{N}(x_1, 1)$$

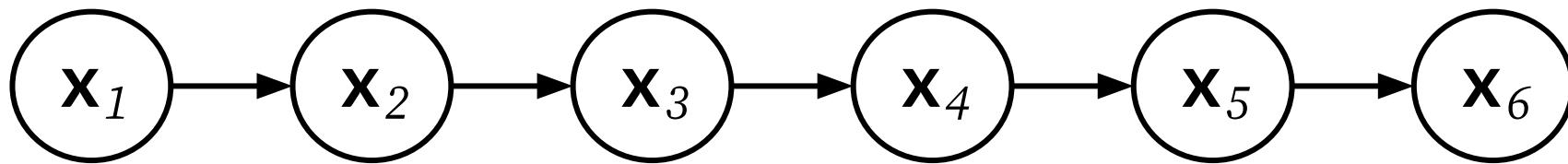
...

$$x_n \sim \mathcal{N}(x_{n-1}, 1)$$

An example (with
normal distributions)

Famous Bayesian Networks – Markov Models

Markov Models ... also known as **Markov chains**.



Generative process:

$$x_1 \sim p(x_1)$$

$$x_2 \sim p(x_2 \mid x_1)$$

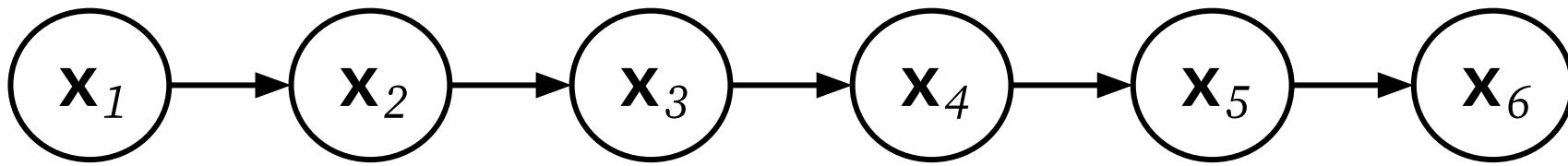
...

$$x_n \sim p(x_n \mid x_{n-1})$$

An alternative way of
writing the generative
process

Famous Bayesian Networks – Markov Models

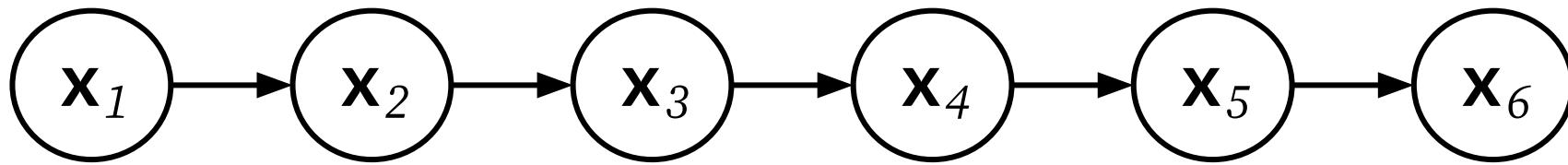
Markov Models ... also known as **Markov chains**.



Applications:

Famous Bayesian Networks – Markov Models

Markov Models ... also known as **Markov chains**.



Applications:

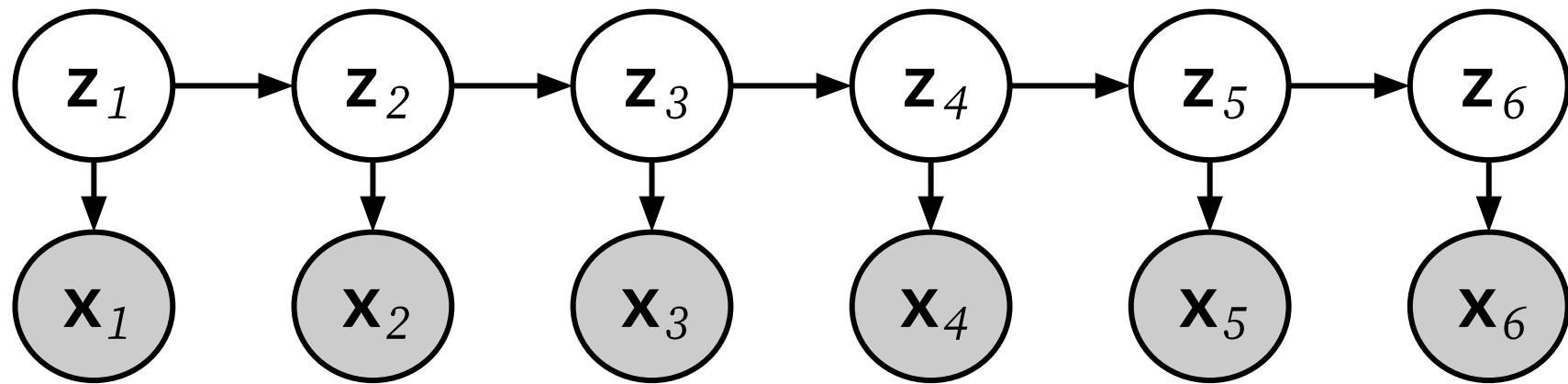
- Markov chain Monte Carlo – an approximate Bayesian inference method.
 - For sampling from a conditional (posterior) distribution of a joint model.
- Business: consumer behavior, brand loyalty, brand switching.
- Operations Research: e.g., models in queuing theory.

Often: sampling, rather than inference of conditional distribution given observations.

Famous Bayesian Networks – Hidden Markov Models (HMMs)

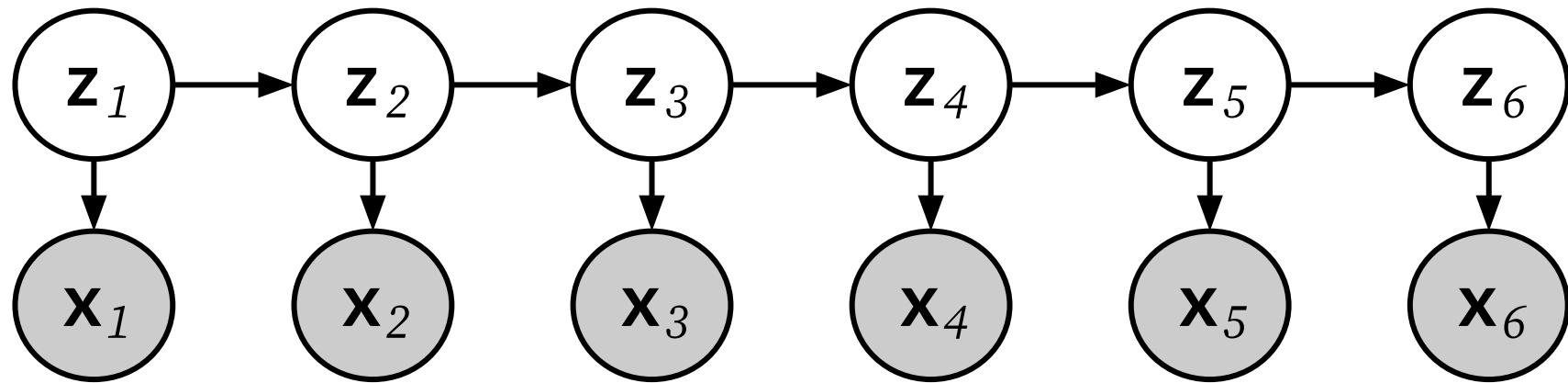
Famous Bayesian Networks – Hidden Markov Models (HMMs)

Hidden Markov Models (HMMs) – an extension to partially observed Markov chains.



Famous Bayesian Networks – Hidden Markov Models (HMMs)

Hidden Markov Models (HMMs) – an extension to partially observed Markov chains.



Backbone of latent variables (in a Markov chain).

One observed variable (“emission variable”) for each latent variable.

Famous Bayesian Networks – Hidden Markov Models (HMMs)

Hidden Markov Models (HMMs) – an extension to partially observed Markov chains.

Generative process:

$$z_1 \sim p(z_1)$$

$$x_1 \sim p(x_1 | z_1)$$

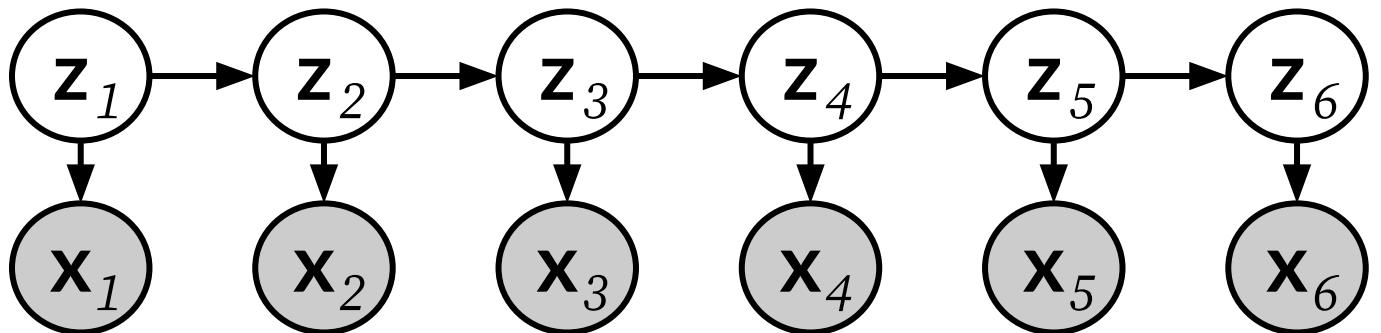
$$z_2 \sim p(z_2 | z_1)$$

$$x_2 \sim p(x_2 | z_2)$$

...

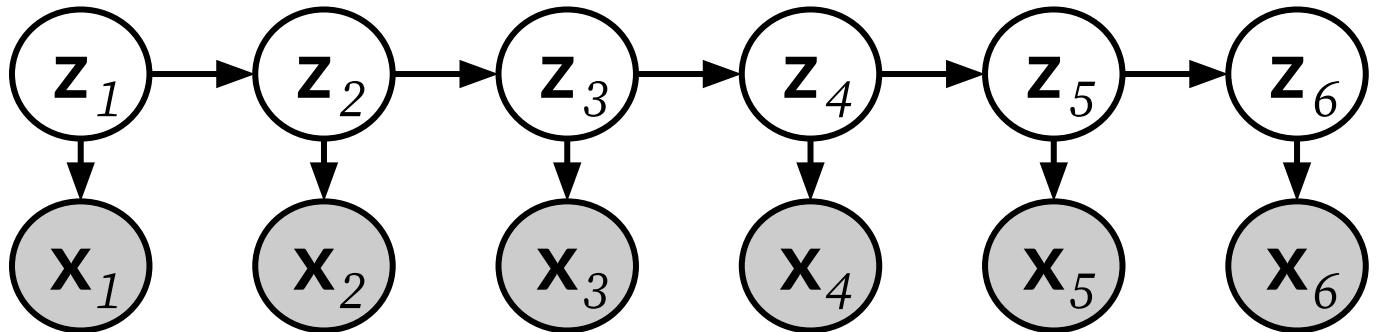
$$z_n \sim p(z_n | z_{n-1})$$

$$x_n \sim p(x_n | z_n)$$



Famous Bayesian Networks – Hidden Markov Models (HMMs)

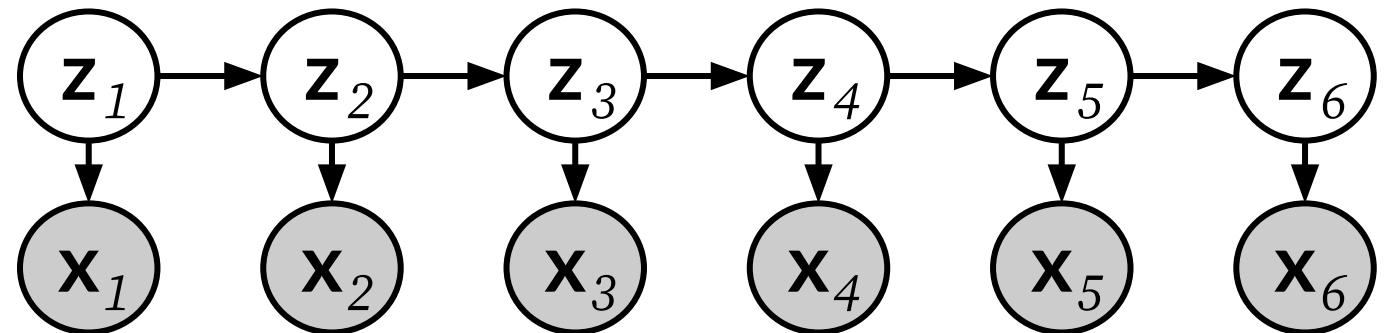
Hidden Markov Models (HMMs) – an extension to partially observed Markov chains.



Applications

Famous Bayesian Networks – Hidden Markov Models (HMMs)

Hidden Markov Models (HMMs) – an extension to partially observed Markov chains.



Applications

- Time series modeling in speech, CV, computational biology (whenever there is a noisy signal).
- Viterbi algorithm (applies to this model): speech recognition, computational linguistics, bioinformatics, digital cellular/satellite/deep-space communications.

Famous Bayesian Networks – Mixture Models (e.g., GMMs)

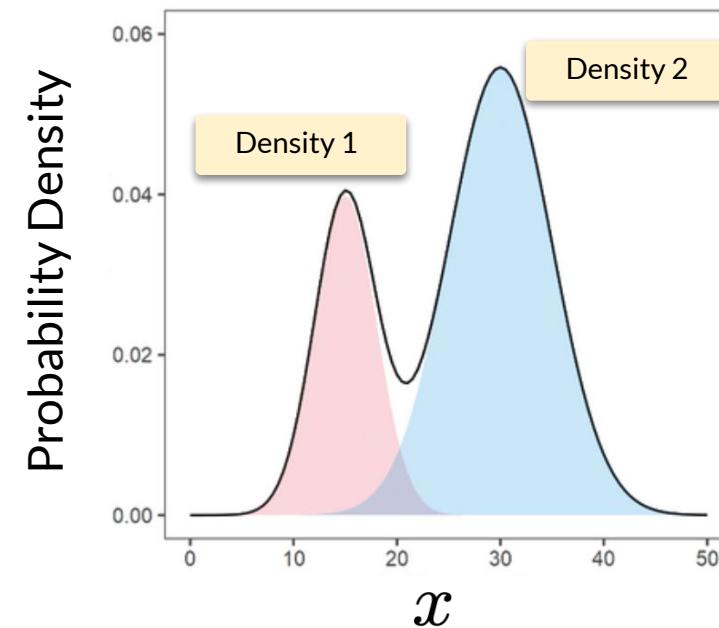
Mixture Models, e.g., Gaussian mixture models (GMMs)

Famous Bayesian Networks – Mixture Models (e.g., GMMs)

Mixture Models, e.g., Gaussian mixture models (GMMs)

Suppose we have samples from a *mixture of distributions* – i.e.

One dimensional
data



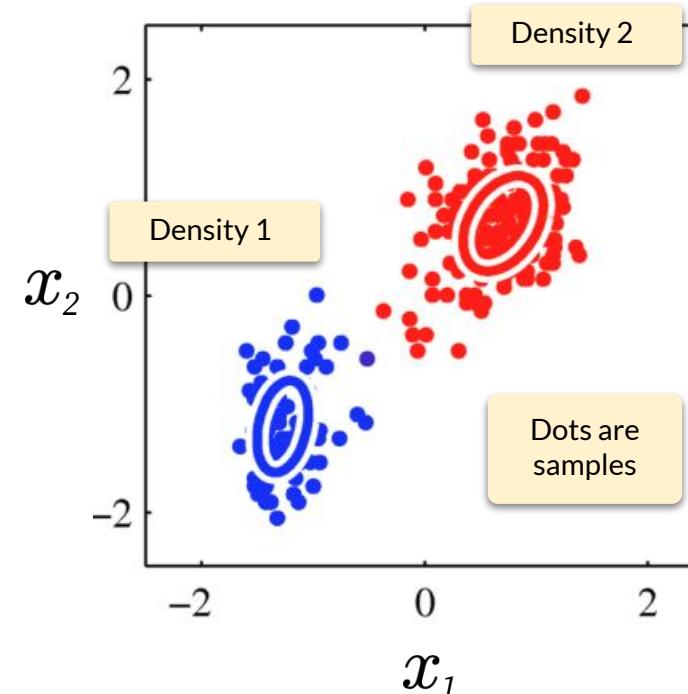
Sometimes called:
“Mixture densities”
or
“Component densities”

Famous Bayesian Networks – Mixture Models (e.g., GMMs)

Mixture Models, e.g., Gaussian mixture models (GMMs)

Suppose we have samples from a *mixture of distributions* – i.e.

Two dimensional
data

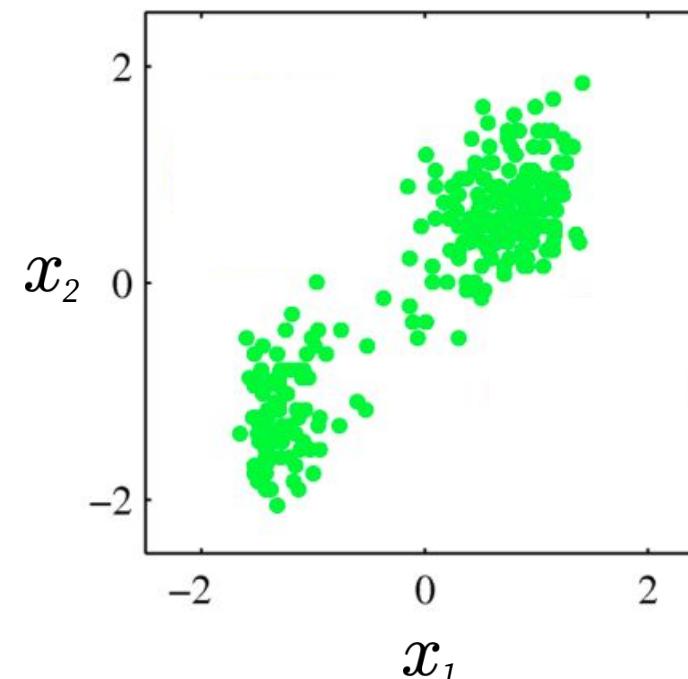


Famous Bayesian Networks – Mixture Models (e.g., GMMs)

Mixture Models, e.g., Gaussian mixture models (GMMs)

Suppose we have samples from a *mixture of distributions* – i.e.

Two dimensional
data



What if we only observe
the samples?

Q: How can we infer the
parameters of the
distributions that
produced them?

Famous Bayesian Networks – Mixture Models (e.g., GMMs)

Mixture Models, e.g., Gaussian mixture models (GMMs)

A: We can use a PGM!

Famous Bayesian Networks – Mixture Models (e.g., GMMs)

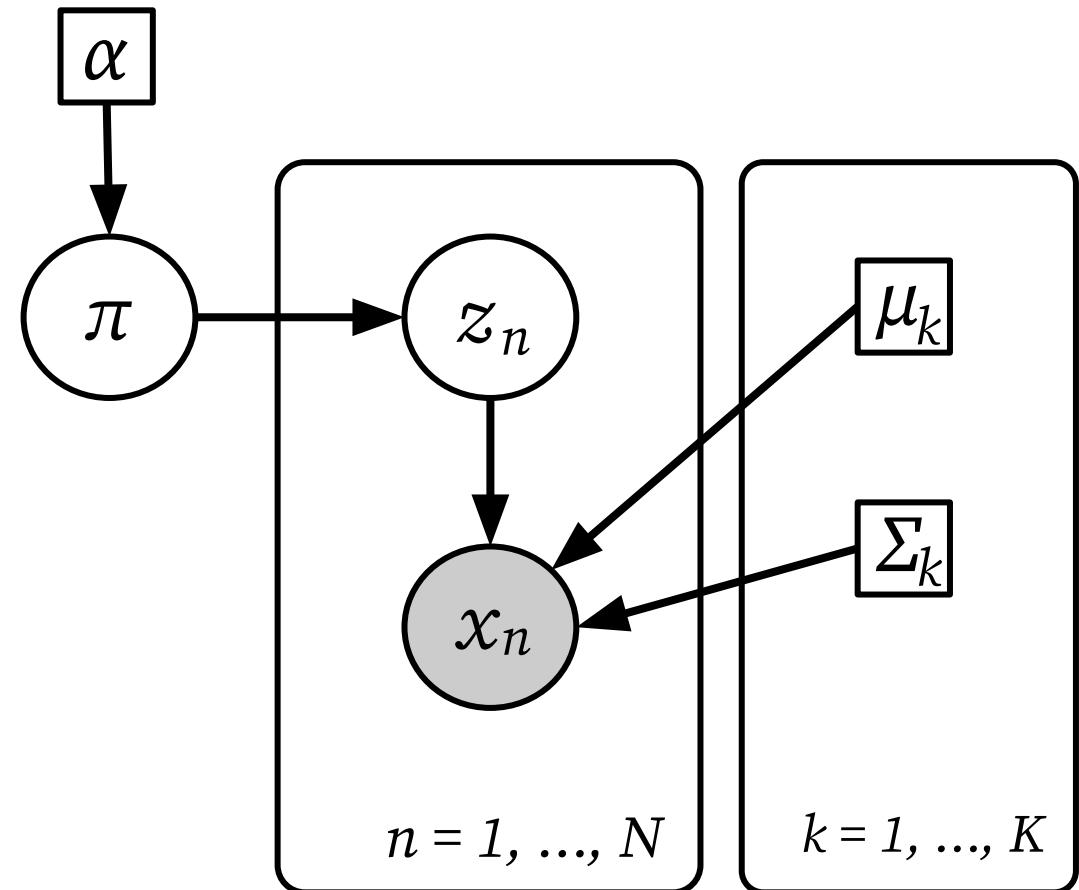
Mixture Models, e.g., Gaussian mixture models (GMMs)

A: We can use a PGM!

Called a **mixture model**.

(Or a Gaussian mixture model, GMM,
if component densities are normal)

Inference/estimation in this
model (given observations)
yields estimate of densities.



Famous Bayesian Networks – Mixture Models (e.g., GMMs)

Mixture Models, e.g., Gaussian mixture models (GMMs)

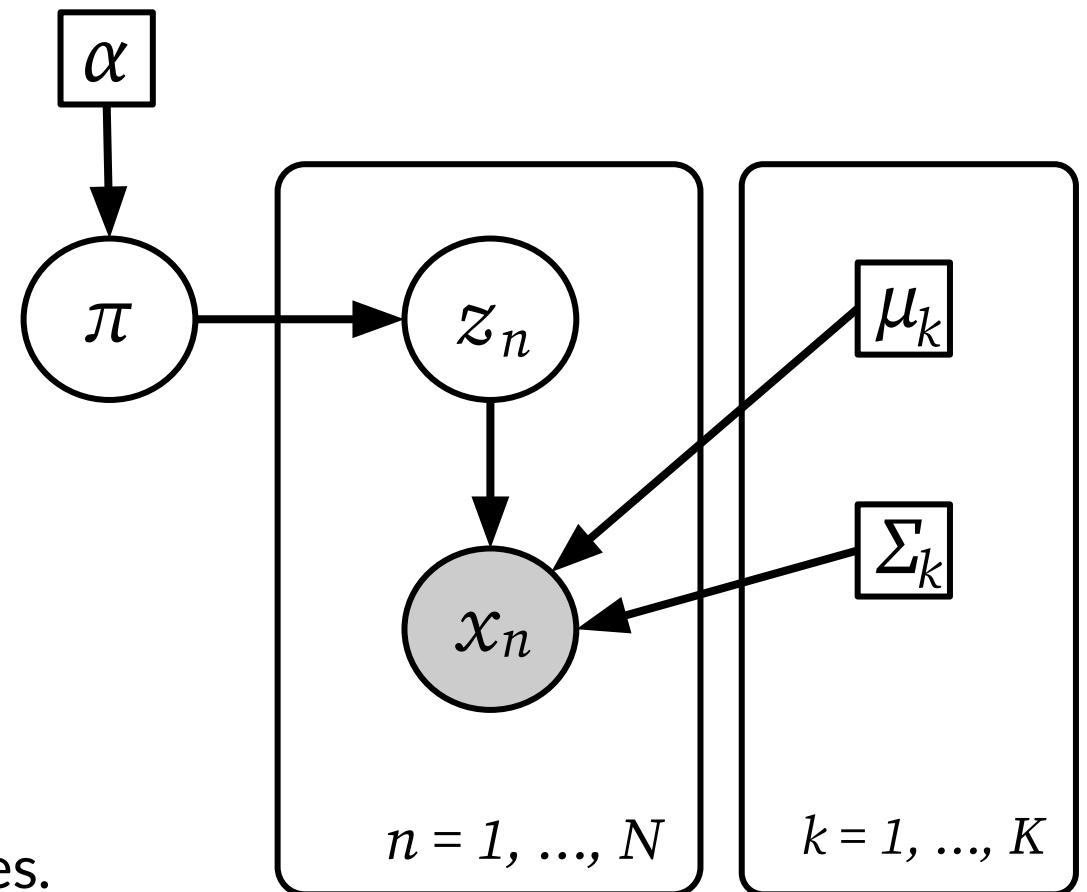
Variables in PGM:

π - Mixture weights

z_n - Assignment variables
(one per observation)

x_n - Observation variables
(corresponds to dots in prev. image)

μ_k, Σ_k - Parameters of mixture densities.



Famous Bayesian Networks – Mixture Models (e.g., GMMs)

Mixture Models, e.g., Gaussian mixture models (GMMs)

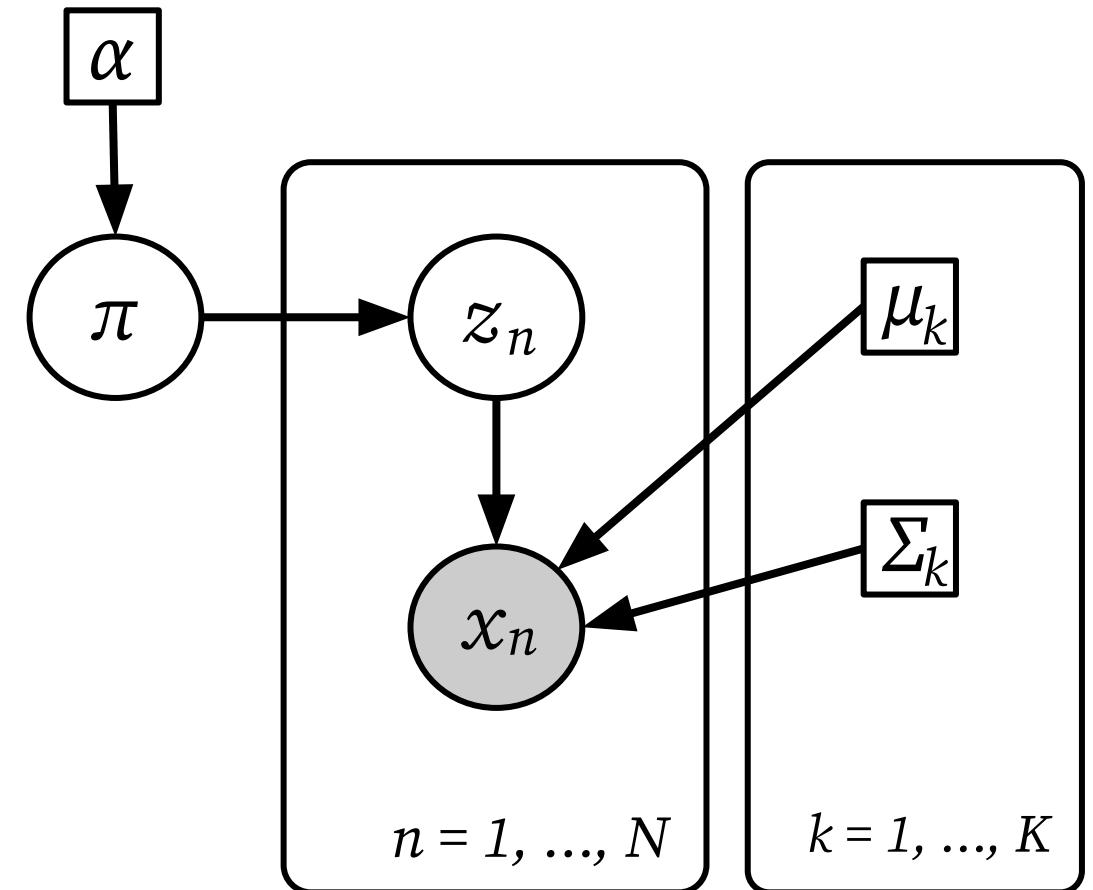
Generative process:

$$\pi \sim \text{Dirichlet}(\alpha)$$

for $n = 1, \dots, N$:

$$z_n \sim \text{Categorical}(\pi)$$

$$x_n \sim \mathcal{N}(\mu_{z_n}, \Sigma_{z_n})$$

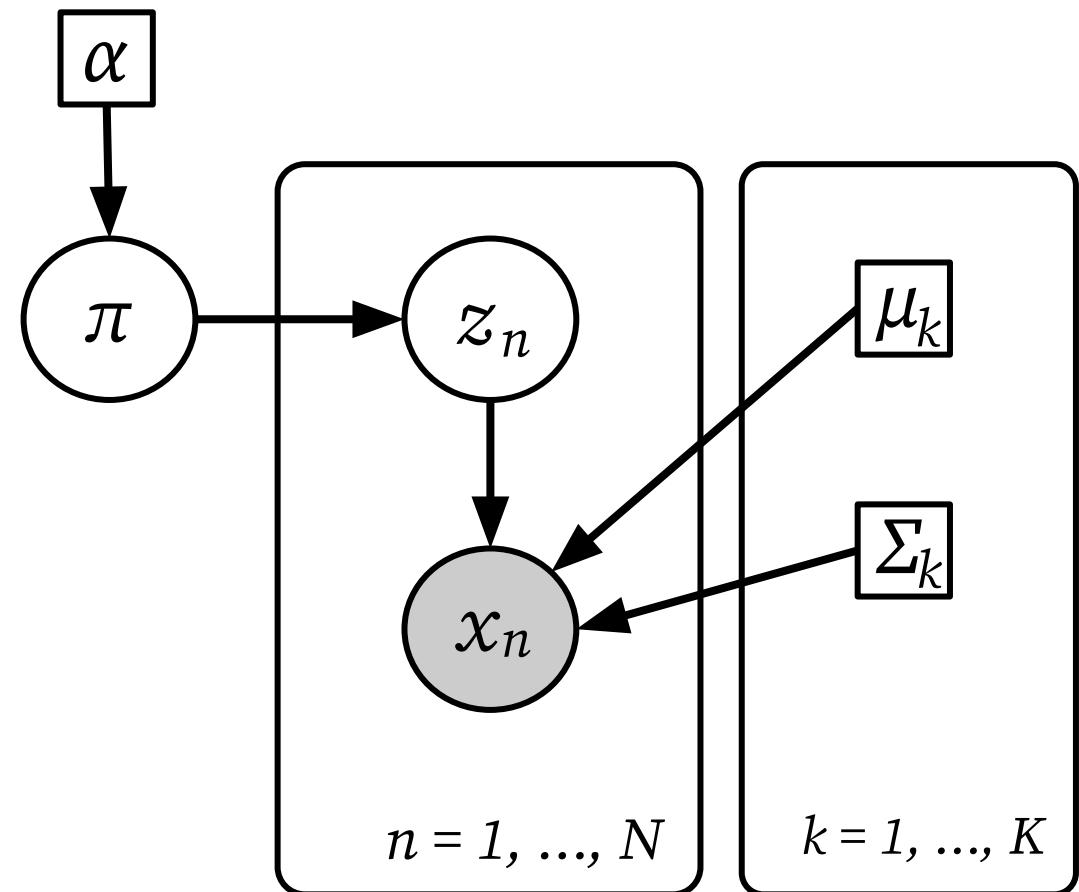


Famous Bayesian Networks – Mixture Models (e.g., GMMs)

Mixture Models, e.g., Gaussian mixture models (GMMs)

Applications:

- Clustering of data.
- Density estimation.
- Anomaly detection.
- Image segmentation.
- Sometimes with far richer mixture densities (using neural networks).



Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Original Paper (2003)
→ 55,776 citations



review articles

Surveying a suite of algorithms that offer a solution to managing large document archives.

BY DAVID M. BLEI

Probabilistic Topic Models

AS OUR COLLECTIVE knowledge continues to be digitized and stored—in the form of news, blogs, Web pages, scientific articles, books, images, sound, video, and social networks—it becomes more difficult to find and discover what we are looking for. We need new computational tools to help organize, search, and understand these vast amounts of information.

Right now, we work with online information using two main tools—search and links. We type keywords into a search engine and find a set of documents related to them. We look at the documents in that set, possibly navigating to other linked documents. This is a powerful way of interacting with our online archive, but something is missing.

Imagine searching and exploring documents based on the themes that run through them. We might “zoom in” and “zoom out” to find specific or broader themes; we might look at how those themes changed through time or how they are connected to each other. Rather than finding documents through keyword search alone, we might first find the theme that we are interested in, and then examine the documents related to that theme.

■ Topic modeling algorithms can be applied to many kinds of data, including other applications they have been used to find patterns in genetic data, images, and social networks.

DOI:10.1145/2133806.2133826

APRIL 2012 | VOL. 55 | NO. 4 | COMMUNICATIONS OF THE ACM 77

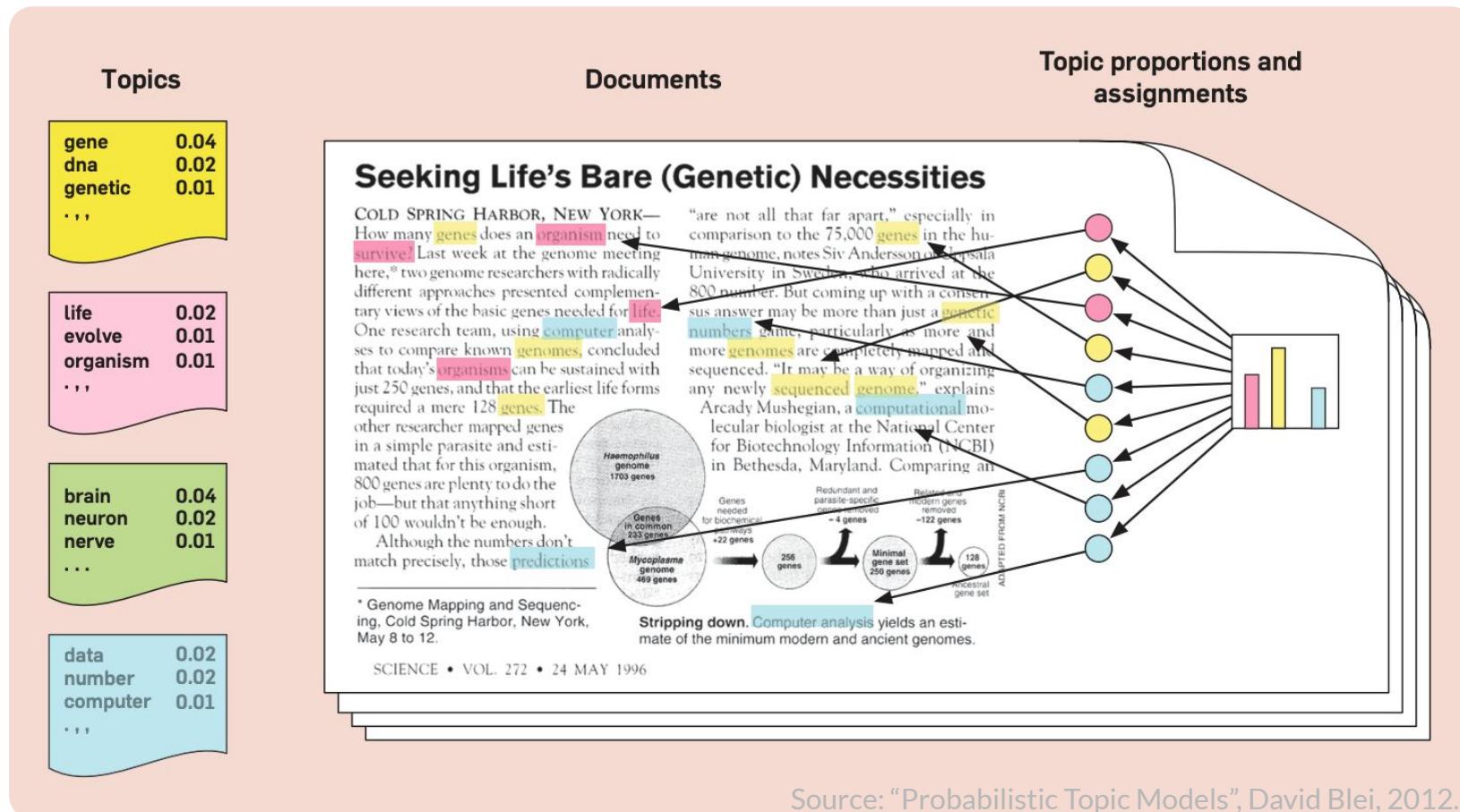
Newer Paper (2012)
“Probabilistic Topic Models”

Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model* over documents.

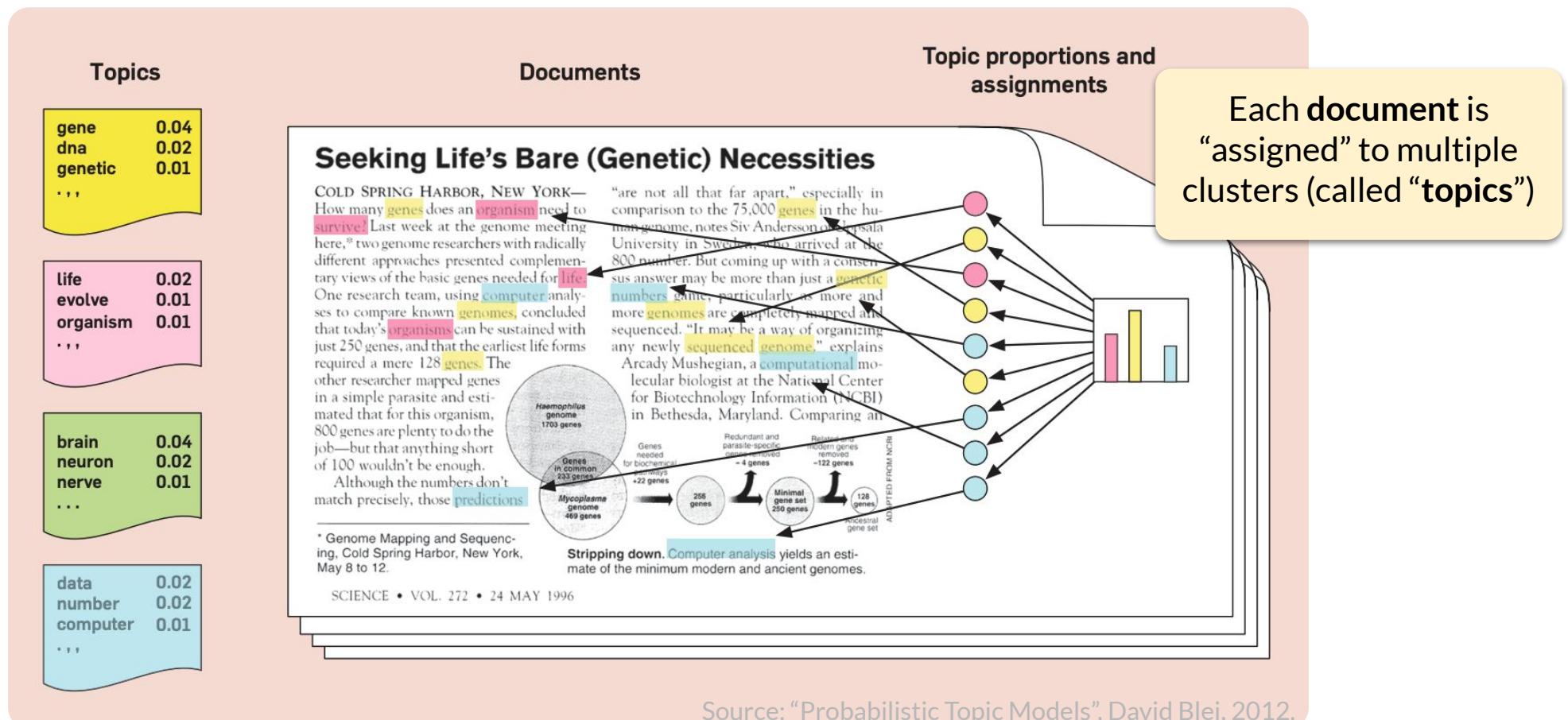
Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model* over documents.



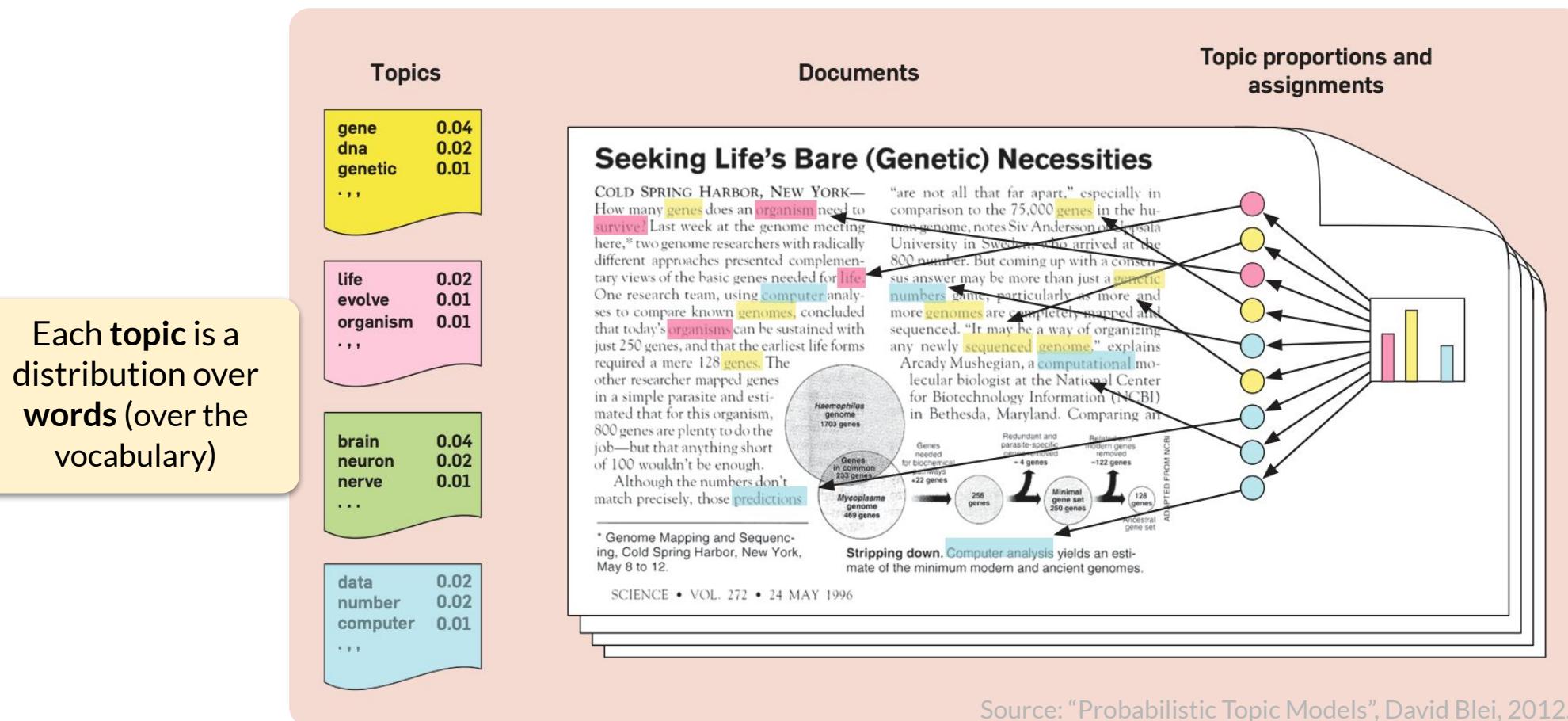
Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model* over documents.



Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

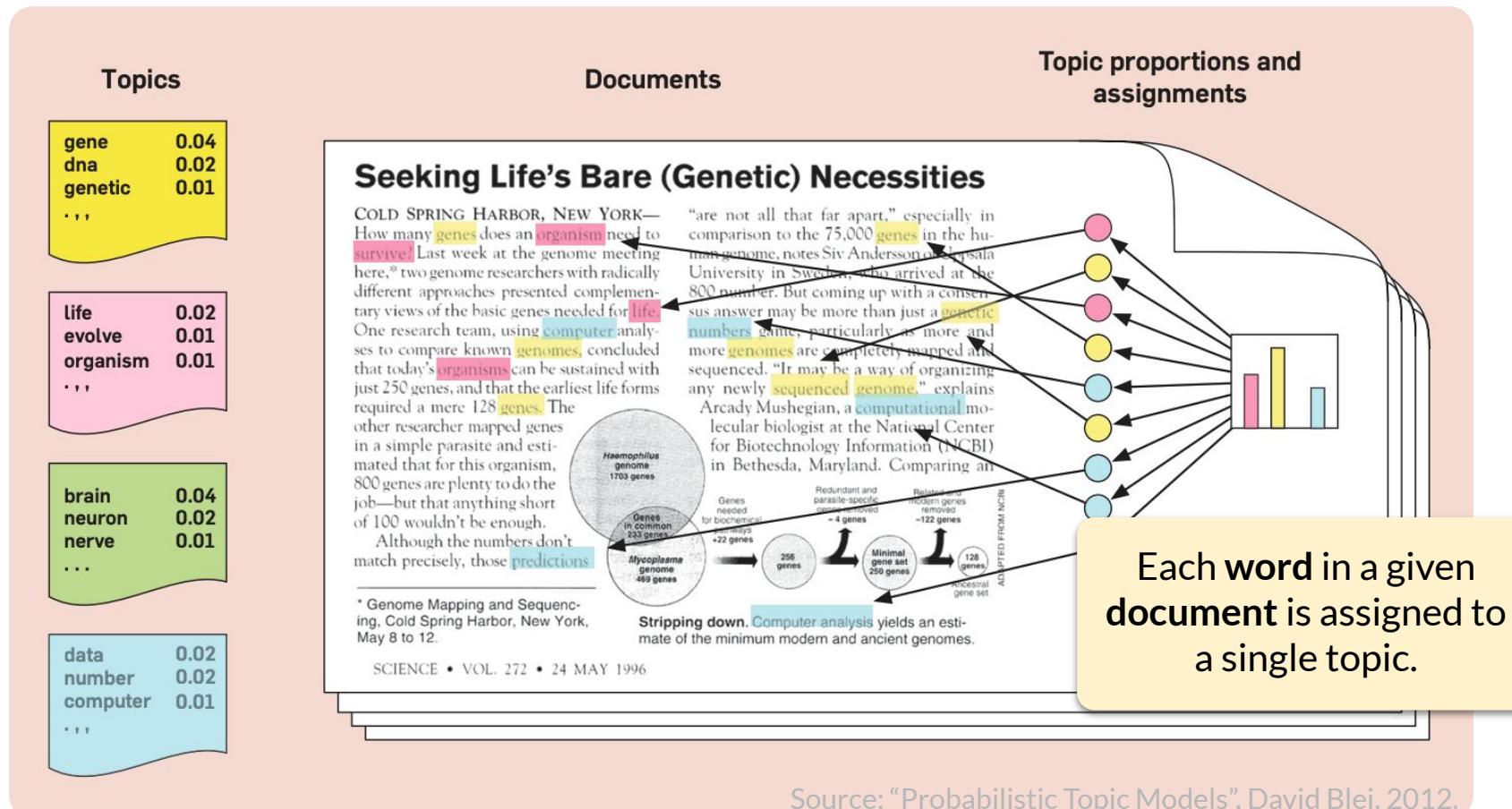
Latent Dirichlet Allocation (LDA) – an *admixture model* over documents.



Source: "Probabilistic Topic Models", David Blei, 2012.

Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

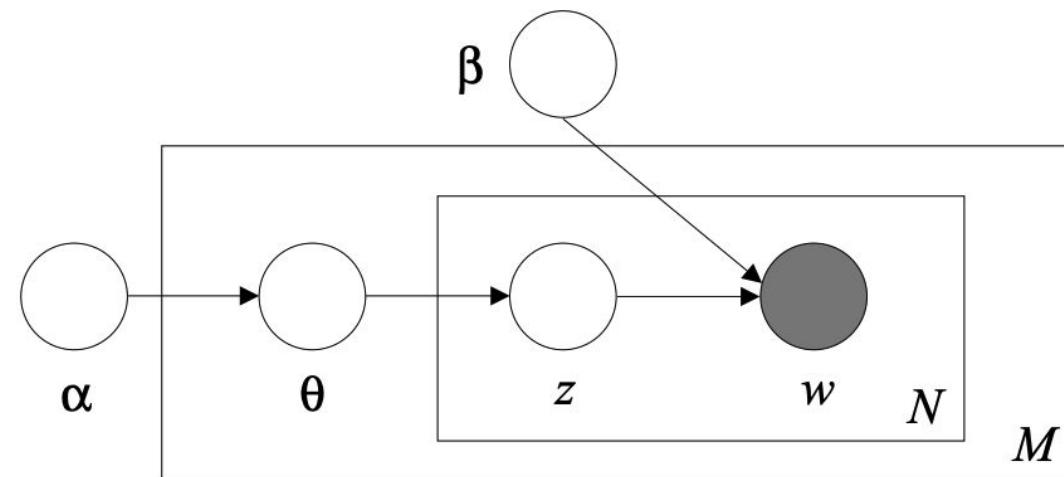
Latent Dirichlet Allocation (LDA) – an *admixture model* over documents.



Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model*.

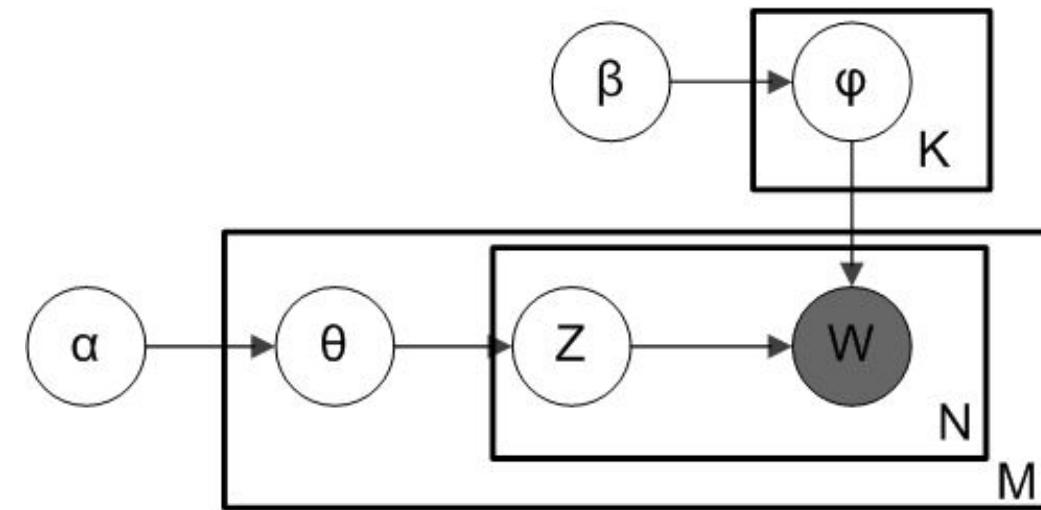
PGM from the original
paper (2003).



Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model*.

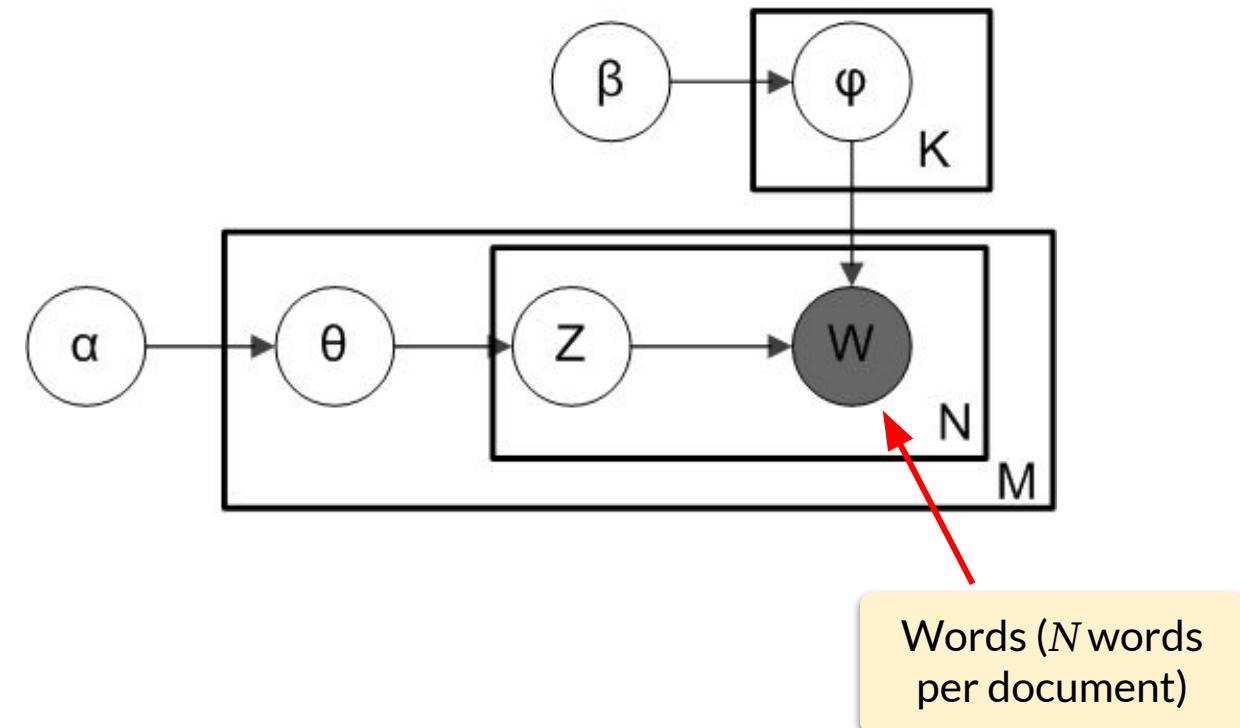
Updated slightly
(from wikipedia)



Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model*.

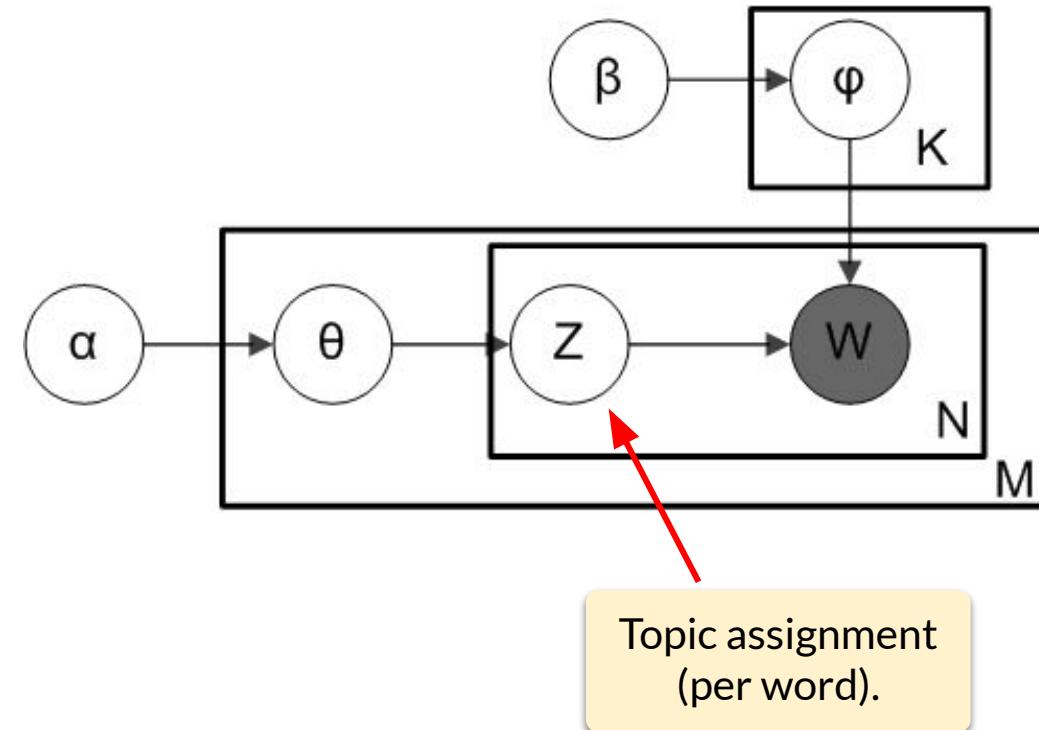
Updated slightly
(from wikipedia)



Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model*.

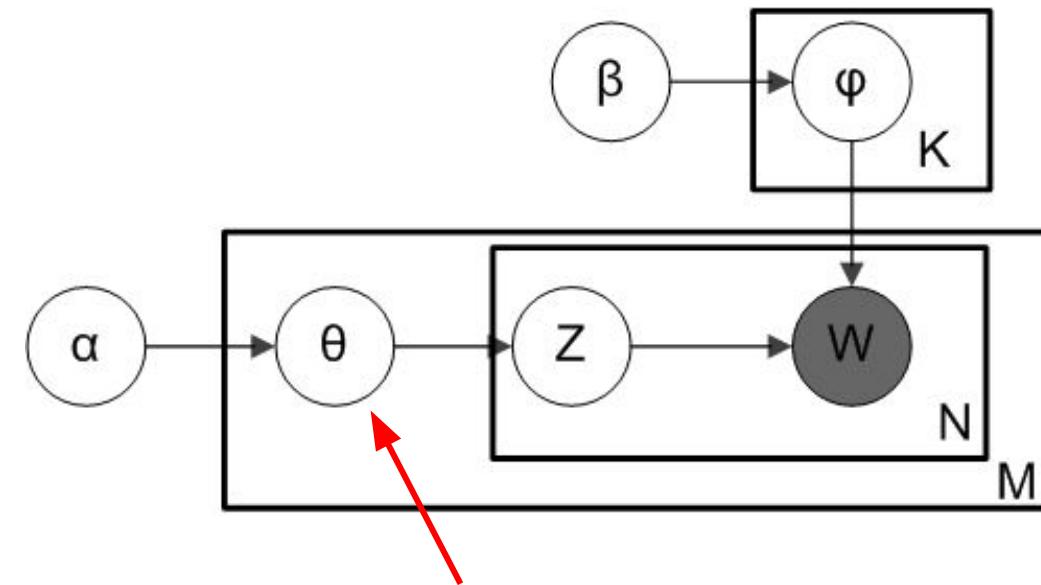
Updated slightly
(from wikipedia)



Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model*.

Updated slightly
(from wikipedia)

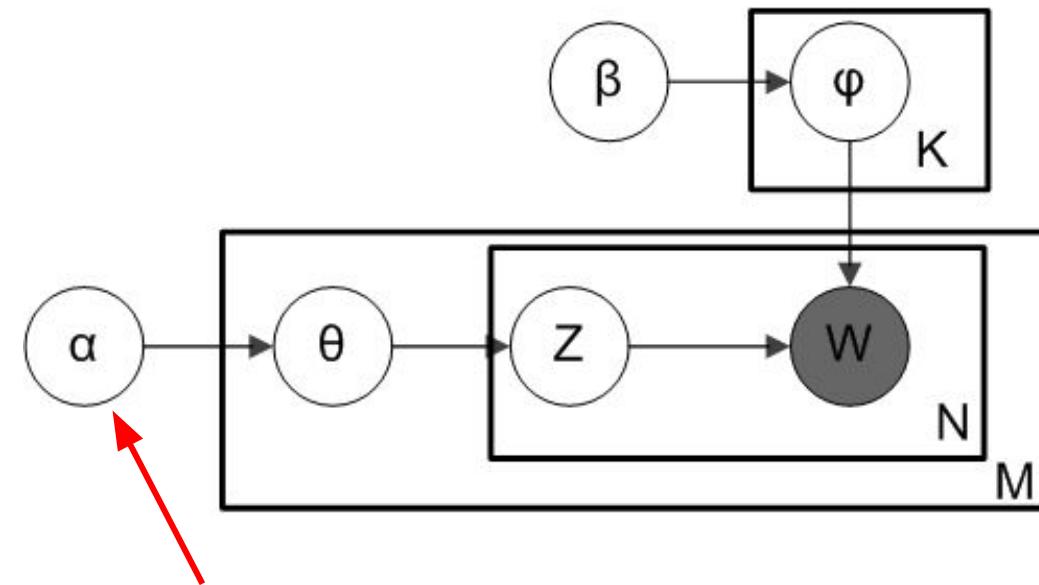


Topic mixture distribution
(for each of M documents)

Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model*.

Updated slightly
(from wikipedia)

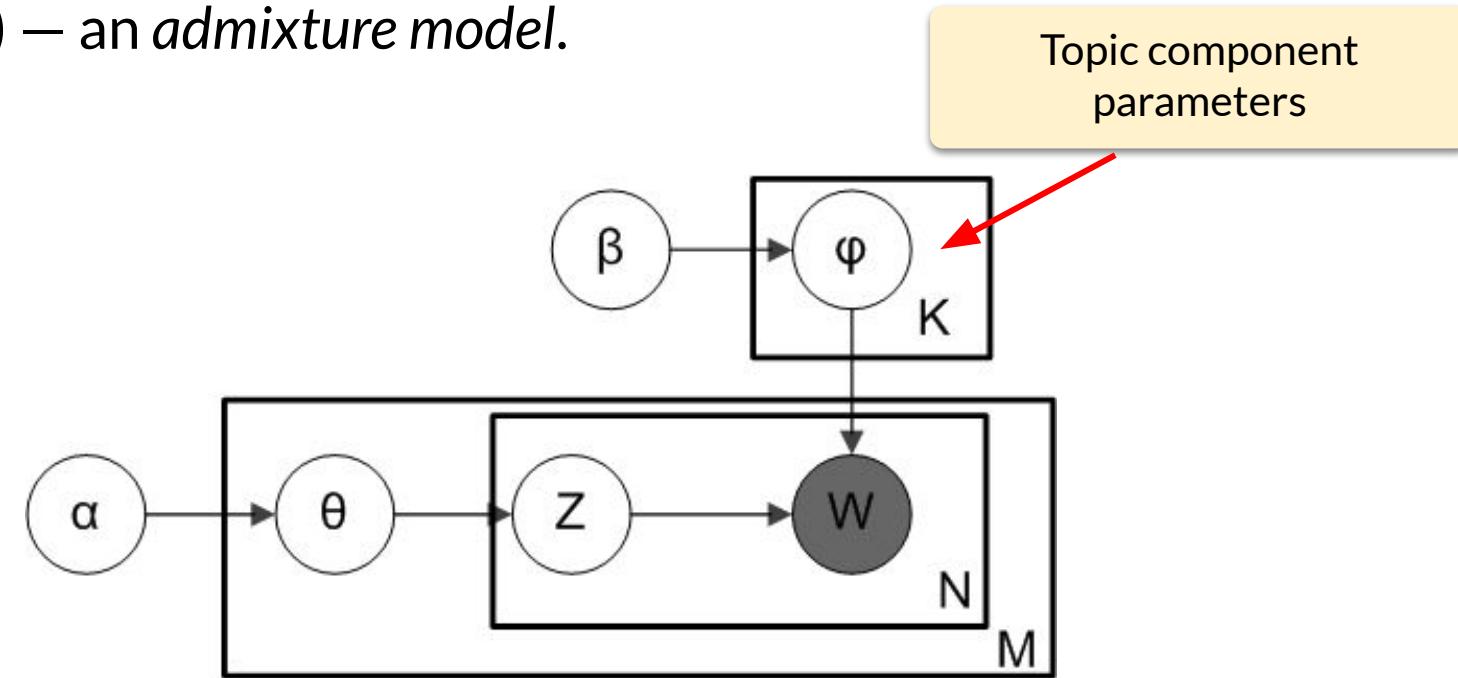


Parameter of prior for
topic distributions.

Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model*.

Updated slightly
(from wikipedia)

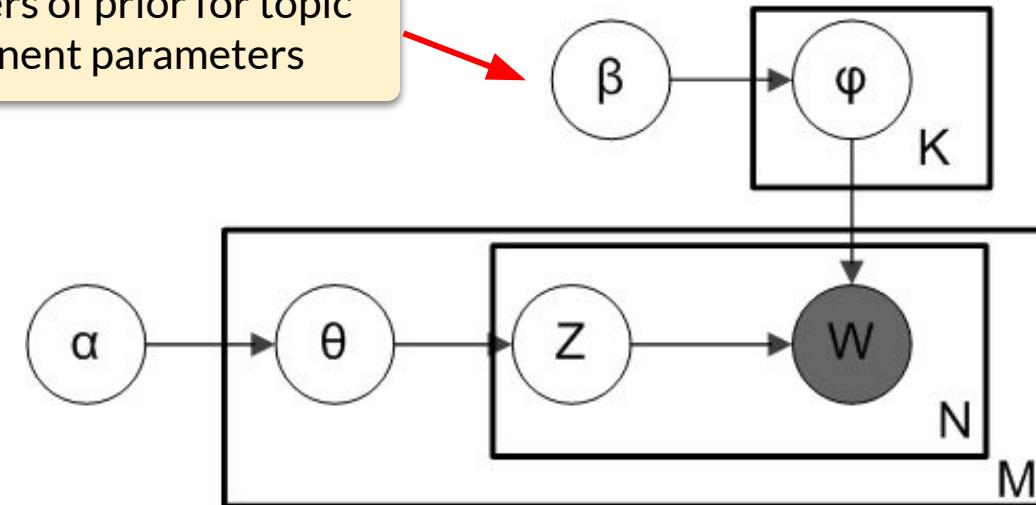


Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model*.

Updated slightly
(from wikipedia)

Parameters of prior for topic
component parameters

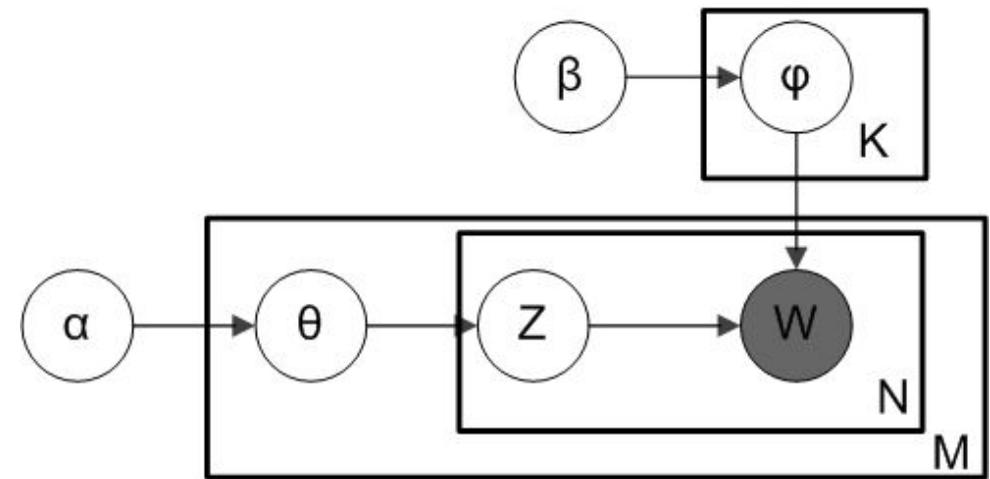


Famous Bayesian Networks – Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) – an *admixture model*.

Applications:

- Information retrieval.
- Content recommendation.
- Recommender system, collaborative filtering (items>users in place of documents>words)
- Document analysis & feature extraction.



Famous Bayesian Networks – Neural Transformations in PGMs

Famous Bayesian Networks – Neural Transformations in PGMs

Can incorporate more-complex function transformations (or neural networks) into Bayesian networks!

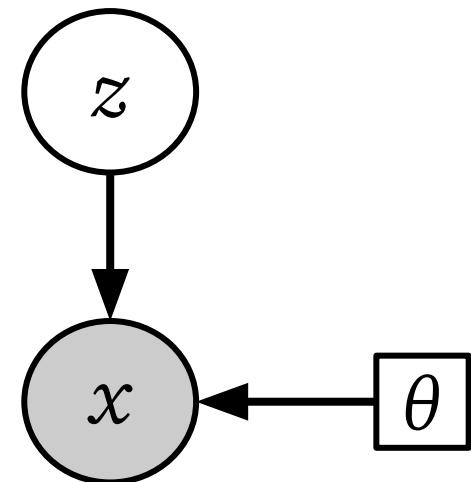
Famous Bayesian Networks – Neural Transformations in PGMs

Can incorporate more-complex function transformations (or neural networks) into Bayesian networks!

For example, a simple PGM might be:

$$z \sim p(z) = \mathcal{N}(0, 1)$$

$$x \sim p(x | z) = \mathcal{N}(z, 1)$$



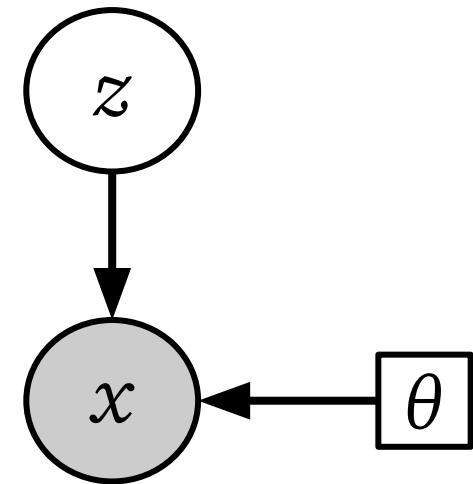
Famous Bayesian Networks – Neural Transformations in PGMs

Can incorporate more-complex function transformations (or neural networks) into Bayesian networks!

But instead we could do:

$$z \sim p(z) = \mathcal{N}(0, 1)$$

$$x \sim p_\theta(x | z) = \mathcal{N}(g_\theta(z), 1)$$



Famous Bayesian Networks – Neural Transformations in PGMs

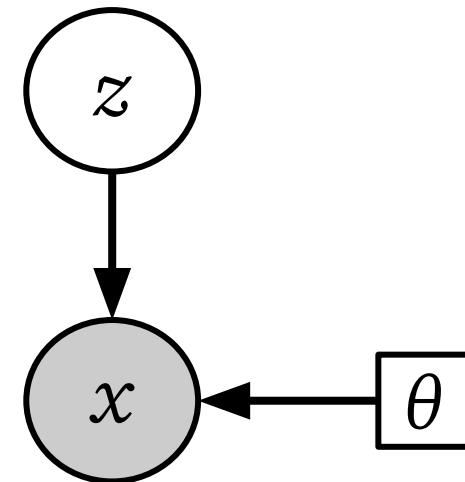
Can incorporate more-complex function transformations (or neural networks) into Bayesian networks!

But instead we could do:

$$z \sim p(z) = \mathcal{N}(0, 1)$$

$$x \sim p_\theta(x | z) = \mathcal{N}(g_\theta(z), 1)$$

$g_\theta(z)$ might be a complex fixed/known, or have parameters that we learn (e.g., if we use a neural network!).



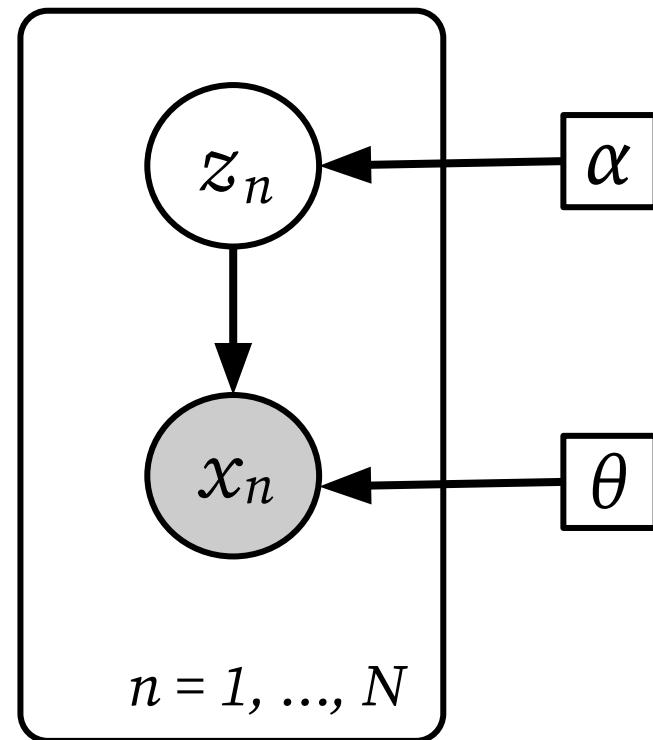
Famous Bayesian Networks – Variational Autoencoder

Example: Variational Autoencoder (VAE), as a Bayesian network.

Famous Bayesian Networks – Variational Autoencoder

Example: Variational Autoencoder (VAE), as a Bayesian network.

Super simple structure!



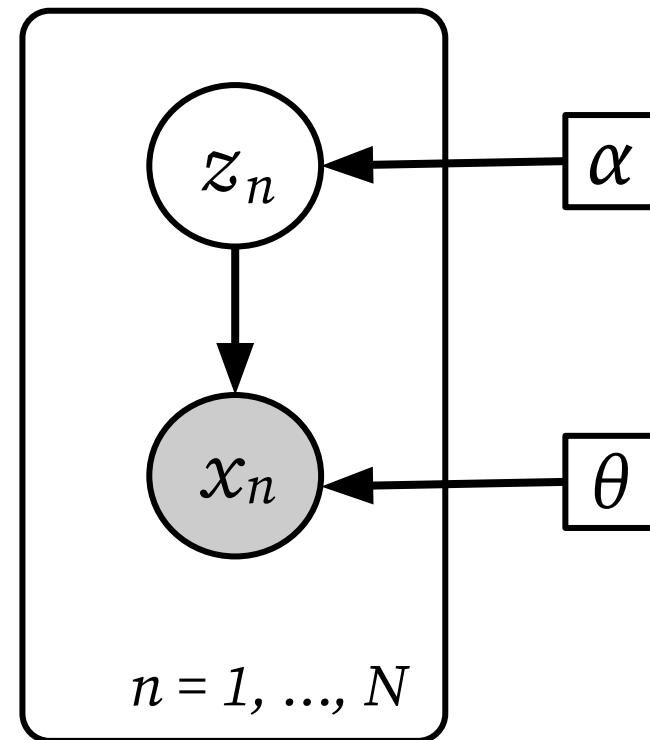
Famous Bayesian Networks – Variational Autoencoder

Example: Variational Autoencoder (VAE), as a Bayesian network.

Super simple structure!

(Will cover VAEs in more detail later in the semester).

Generative process defined by a neural network – with parameters θ .



Famous Bayesian Networks – Variational Autoencoder

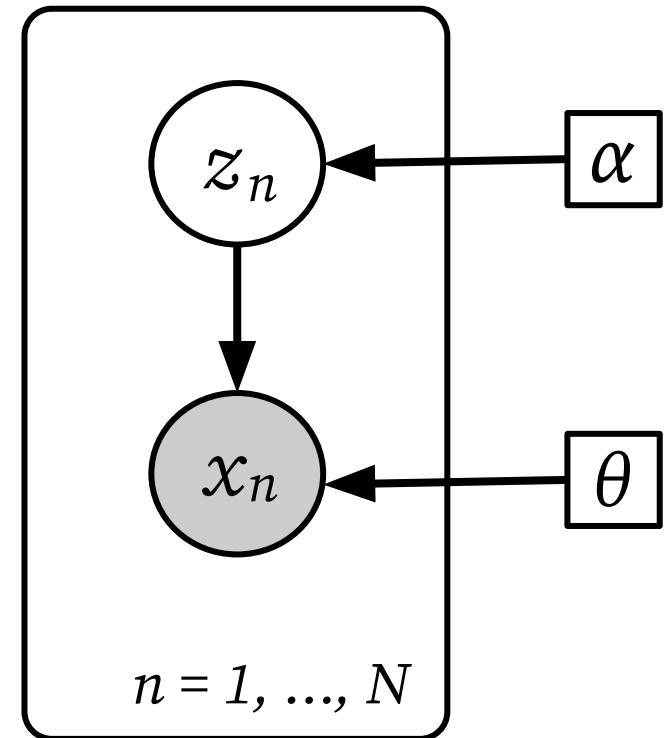
Example: Variational Autoencoder (VAE), as a Bayesian network.

Generative process:

for $n = 1, \dots, N$:

$$z_n \sim p(z | \alpha)$$

$$x_n \sim p_\theta(x | z_n)$$



Famous Bayesian Networks – Variational Autoencoder

Example: Variational Autoencoder (VAE), as a Bayesian network.

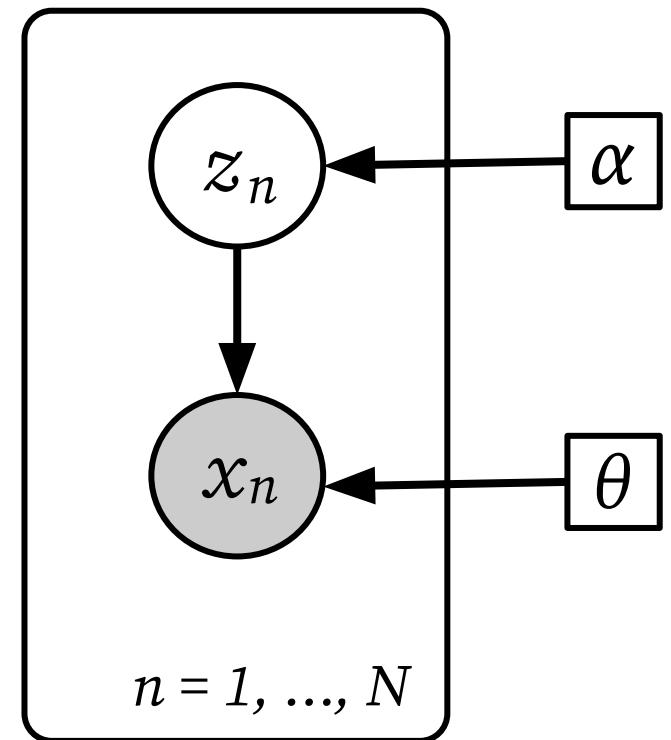
Generative process:

For each data point n

for $n = 1, \dots, N$:

$$z_n \sim p(z | \alpha)$$

$$x_n \sim p_\theta(x | z_n)$$



Famous Bayesian Networks – Variational Autoencoder

Example: Variational Autoencoder (VAE), as a Bayesian network.

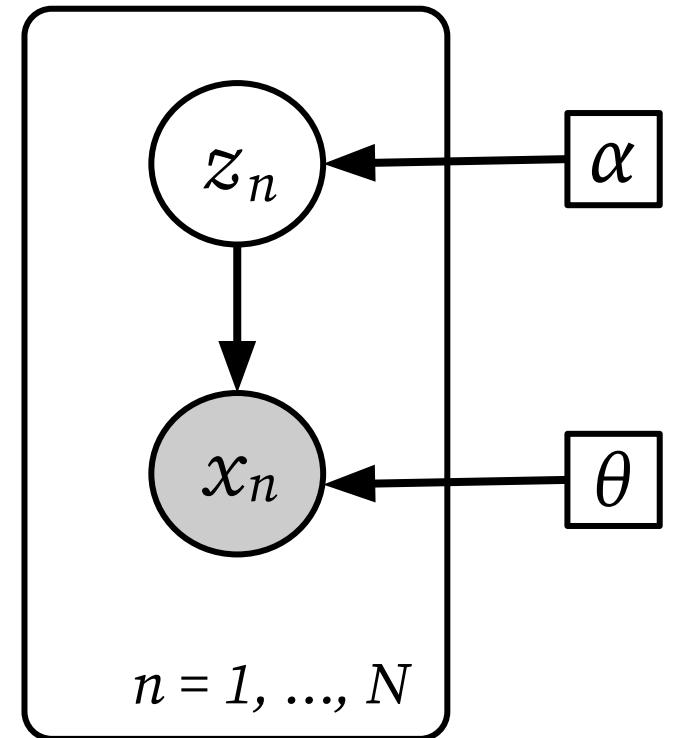
Generative process:

for $n = 1, \dots, N$:

Sample from prior

$$z_n \sim p(z | \alpha)$$

$$x_n \sim p_\theta(x | z_n)$$



Famous Bayesian Networks – Variational Autoencoder

Example: Variational Autoencoder (VAE), as a Bayesian network.

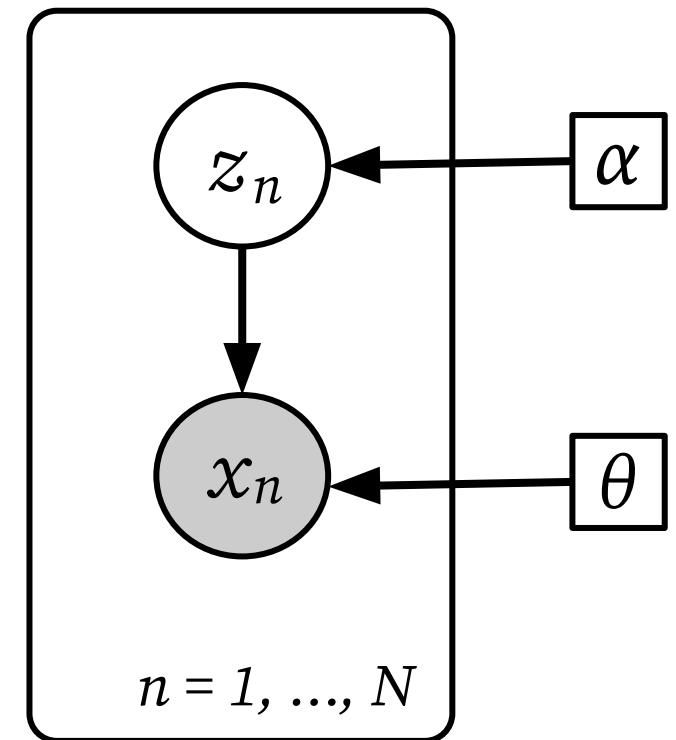
Generative process:

for $n = 1, \dots, N$:

$$z_n \sim p(z | \alpha)$$

Sample from distribution
parameterized via θ .

$$x_n \sim p_\theta(x | z_n)$$



Famous Bayesian Networks – Variational Autoencoder

Example: Variational Autoencoder (VAE), as a Bayesian network.

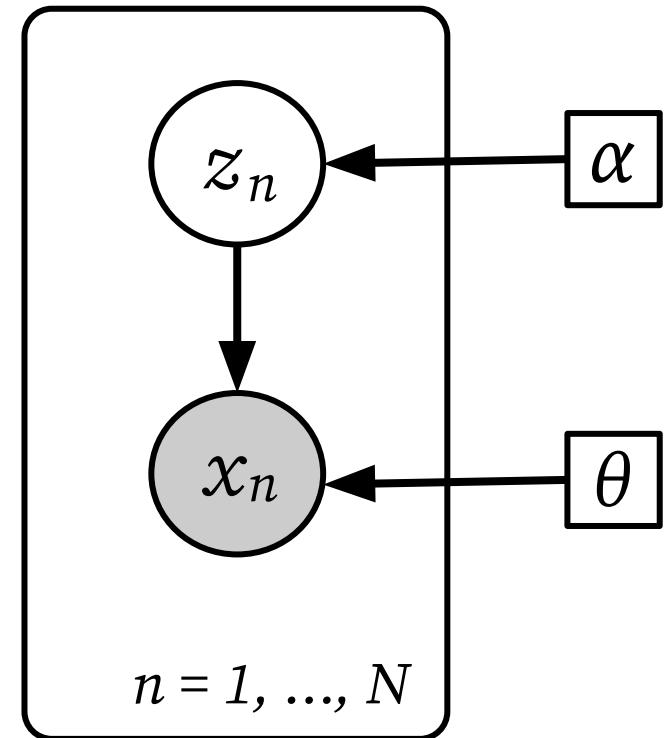
Generative process:

for $n = 1, \dots, N$:

$$z_n \sim p(z | \alpha) = \mathcal{N}(0, \alpha^2)$$

$$x_n \sim p_\theta(x | z_n) = \mathcal{N}(\mu_\theta(z_n), \Sigma_\theta(z_n))$$

Example densities



Famous Bayesian Networks – Variational Autoencoder

Example: Variational Autoencoder (VAE), as a Bayesian network.

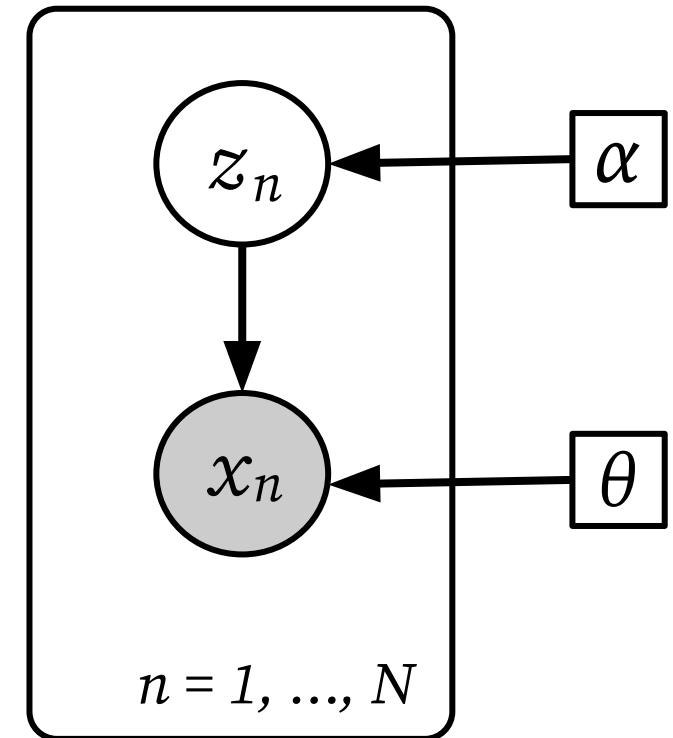
Generative process:

for $n = 1, \dots, N$:

$$z_n \sim p(z | \alpha) = \mathcal{N}(0, \alpha^2)$$

$$x_n \sim p_\theta(x | z_n) = \mathcal{N}(\mu_\theta(z_n), \Sigma_\theta(z_n))$$

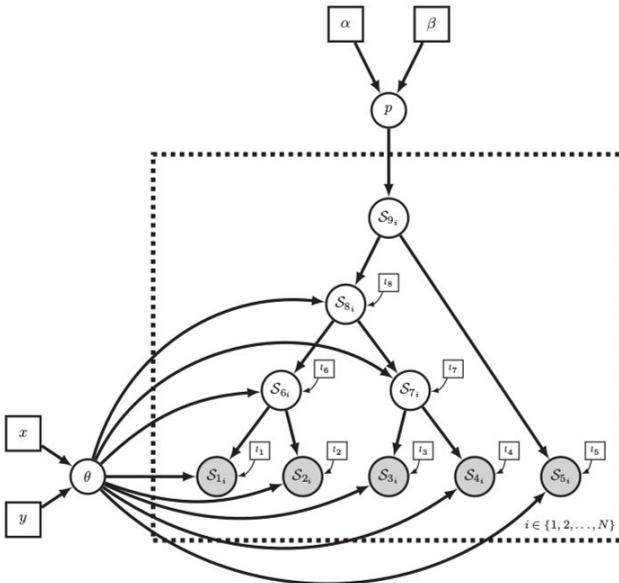
Sample from distribution
parameterized via output of $\mu_\theta(z_n), \Sigma_\theta(z_n)$
a neural network.



Bayesian Networks – Many Other Fun Examples

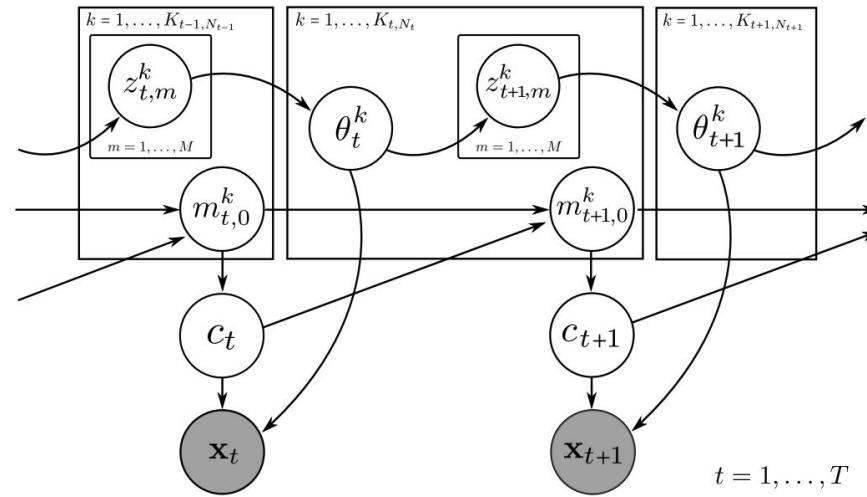
Bayesian Networks – Many Other Fun Examples

Phylogenetics



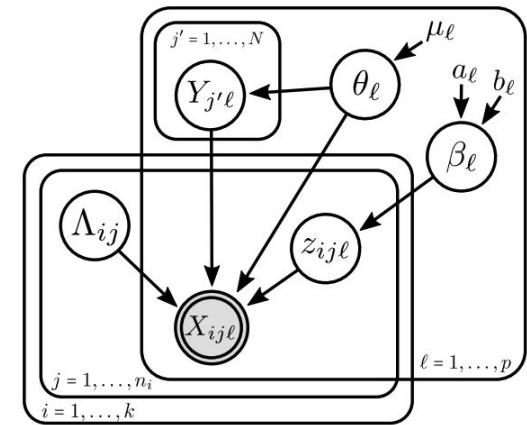
Source: "Probabilistic Graphical Model Representation in Phylogenetics", Höhna 2014

Object Tracking (in videos)



Source: "The Dependent Dirichlet Process Mixture of Objects for Detection-free Tracking and Object Modeling", Neiswanger et al., 2014

Record Linkage (in databases)



Source: "Performance Bounds for Graphical Record Linkage", Steorts et al., 2014

Bayesian Networks – Summary

Bayesian Networks – Summary

Bayesian network is defined by a tuple $(G = \{V, E\}, P)$ where P is specified as a set of local conditional probability distributions (CPDs) associated with nodes V

⇒ Efficient representation using a graph-based data structure.

Computing joint probability of any assignment obtained by multiplying CPDs.

Can sample from joint by sampling from CPDs according to DAG ordering.

Can identify some conditional independencies by looking at graph properties.

Nodes in graph can correspond to random variables or random vectors.

Undirected PGMs: Markov Random Fields

Undirected PGMs – Introduction

Undirected PGMs – Introduction

Bayesian network can represent a rich set of probability distributions (and conditional independence assumptions).

However, some distributions have independence assumptions that cannot be captured exactly.

Undirected PGMs – Introduction

Bayesian network can represent a rich set of probability distributions (and conditional independence assumptions).

However, some distributions have independence assumptions that cannot be captured exactly.

Example. Consider the generative process:

$$X, Y \sim \text{Bernoulli}(0.5)$$

$$Z = X \text{ xor } Y$$

Undirected PGMs – Introduction

Bayesian network can represent a rich set of probability distributions (and conditional independence assumptions).

However, some distributions have independence assumptions that cannot be captured exactly.

Example. Consider the generative process:

$$X, Y \sim \text{Bernoulli}(0.5)$$

$$Z = X \text{ xor } Y$$

1 if X and Y are different.
0 otherwise.

Undirected PGMs – Introduction

Bayesian network can represent a rich set of probability distributions (and conditional independence assumptions).

However, some distributions have independence assumptions that cannot be captured exactly.

Example. Consider the generative process:

$$X, Y \sim \text{Bernoulli}(0.5)$$

$$Z = X \text{ xor } Y$$

$$\Rightarrow \{X \perp Y, \quad X \perp Z, \quad Y \perp Z\}$$

Undirected PGMs – Introduction

Bayesian network can represent a rich set of probability distributions (and conditional independence assumptions).

However, some distributions have independence assumptions that cannot be captured exactly.

Example. Consider the generative process:

$$X, Y \sim \text{Bernoulli}(0.5)$$

$$Z = X \text{ xor } Y$$

$$\Rightarrow \{X \perp Y, \quad X \perp Z, \quad Y \perp Z\}$$

But Z is not independent of $\{X, Y\}$.

Undirected PGMs – Introduction

Bayesian network can represent a rich set of probability distributions (and conditional independence assumptions).

However, some distributions have independence assumptions that cannot be captured exactly.

Example. Consider the generative process:

$$X, Y \sim \text{Bernoulli}(0.5)$$

$$Z = X \text{ xor } Y$$

$$\Rightarrow \{X \perp Y, \quad X \perp Z, \quad Y \perp Z\}$$

But Z is not independent of $\{X, Y\}$.

Called “noisy xor” example.

Undirected PGMs – Introduction

Bayesian network can represent a rich set of probability distributions (and conditional independence assumptions).

However, some distributions have independence assumptions that cannot be captured exactly.

Example. Consider the generative process:

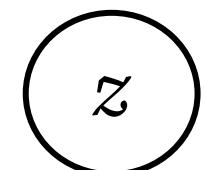
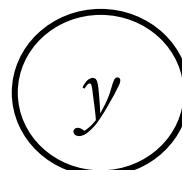
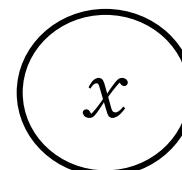
$$X, Y \sim \text{Bernoulli}(0.5)$$

$$Z = X \text{ xor } Y$$

$$\Rightarrow \{X \perp Y, \quad X \perp Z, \quad Y \perp Z\}$$

But Z is not independent of $\{X, Y\}$.

What graph structure to use?



?

Undirected PGMs – Introduction

Bayesian network can represent a rich set of probability distributions (and conditional independence assumptions).

However, some distributions have independence assumptions that cannot be captured exactly.

Example. Consider the generative process:

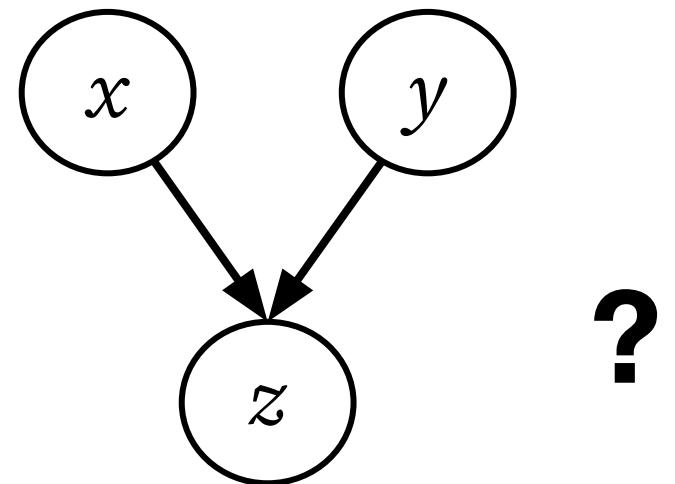
$$X, Y \sim \text{Bernoulli}(0.5)$$

$$Z = X \text{ xor } Y$$

$$\Rightarrow \{X \perp Y, \quad X \perp Z, \quad Y \perp Z\}$$

But Z is not independent of $\{X, Y\}$.

What graph structure to use?



Undirected PGMs – Introduction

Bayesian network can represent a rich set of probability distributions (and conditional independence assumptions).

However, some distributions have independence assumptions that cannot be captured exactly.

Example. Consider the generative process:

$$X, Y \sim \text{Bernoulli}(0.5)$$

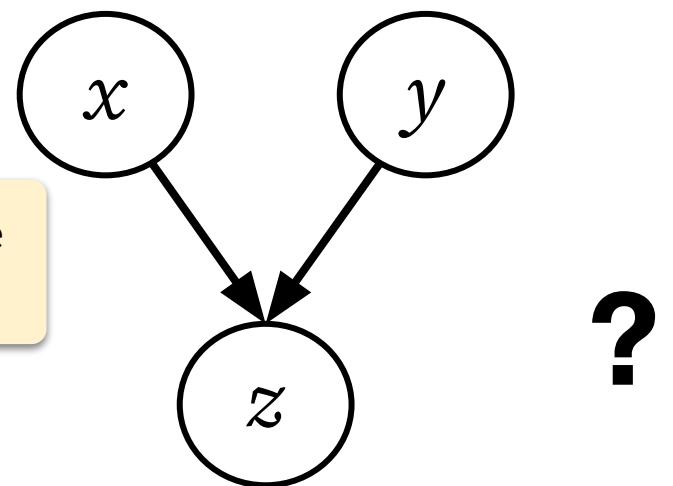
$$Z = X \text{ xor } Y$$

$$\Rightarrow \{X \perp Y, \quad X \perp Z, \quad Y \perp Z\}$$

But Z is **not** independent of $\{X, Y\}$.

PGM does **not** capture all independencies!

What graph structure to use?



Undirected PGMs – Introduction

Bayesian network can represent a rich set of probability distributions (and conditional independence assumptions).

However, some distributions have independence assumptions that cannot be captured exactly.

Furthermore, some probability distributions (*i.e.*, models of the real world) are more easy to define using alternative abstractions...

Undirected PGMs – Introduction

There exists an alternative “**visual language**” for compactly representing and visualizing a probability distribution that is based on **undirected graphs**.

Undirected PGMs – Introduction

There exists an alternative “**visual language**” for compactly representing and visualizing a probability distribution that is based on **undirected graphs**.

These are known as Markov Random Fields (MRFs).

In some cases, they can capture conditional independence assumptions that Bayesian networks cannot represent.

(But come with their own set of drawbacks as well)

Markov Random Fields – Motivating Example

Markov Random Fields – Motivating Example

As a motivating example...

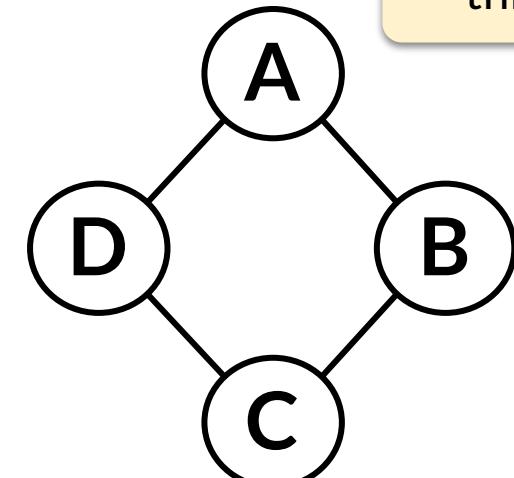
Suppose that we are modeling voting preferences among people: A , B , C , D

Markov Random Fields – Motivating Example

As a motivating example...

Suppose that we are modeling voting preferences among people: A, B, C, D

Let's say that $(A, B), (B, C), (C, D)$, and (D, A) are friends.



Naturally represented by this graph.

Markov Random Fields – Motivating Example

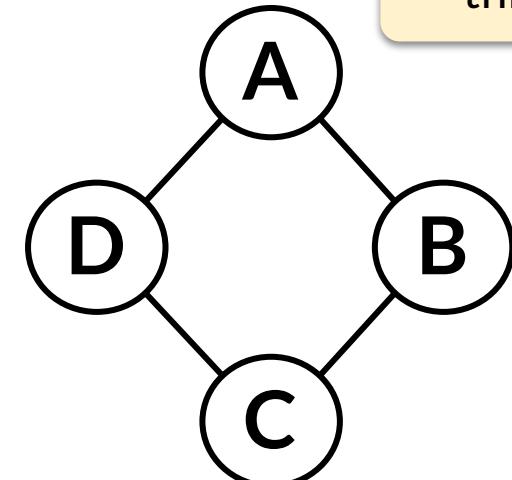
As a motivating example...

Suppose that we are modeling voting preferences among people: A, B, C, D

Let's say that $(A, B), (B, C), (C, D)$, and (D, A) are friends.

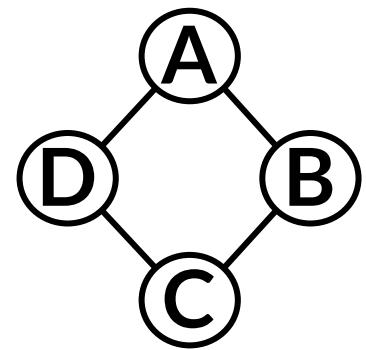
And friends tend to have **similar voting preferences**.

Naturally represented by this graph.



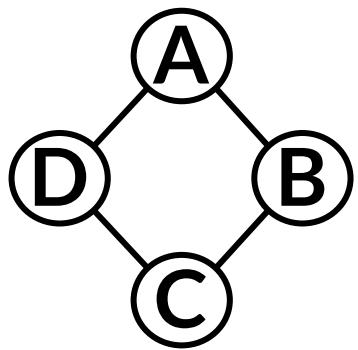
Source: Stefano Ermon, Deep Generative Models (CS236) Class

Markov Random Fields – Motivating Example



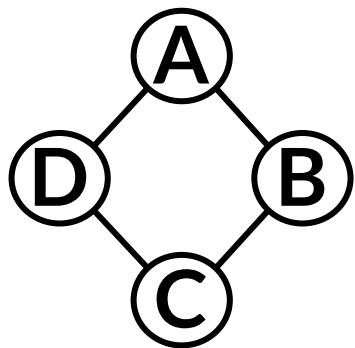
One way to define a probability over the joint voting decision of A, B, C, D is:

Markov Random Fields – Motivating Example



One way to define a probability over the joint voting decision of A, B, C, D is:

- Assign scores to each assignment of these variables, then define probability as the normalized score.



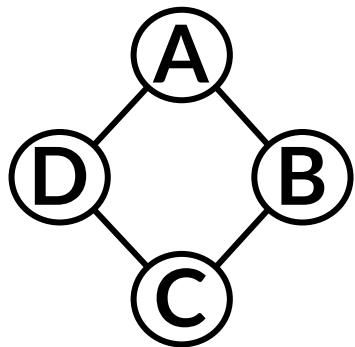
Markov Random Fields – Motivating Example

One way to define a probability over the joint voting decision of A, B, C, D is:

- Assign scores to each assignment of these variables, then define probability as the normalized score.
- A score can be any function, but we'll define it to be of the form:

$$\tilde{p}(A, B, C, D) = \phi(A, B) \phi(B, C) \phi(C, D) \phi(D, A)$$

- Where $\phi(X, Y)$ is a *factor* that assigns more weight to consistent vote among friends X, Y .



Markov Random Fields – Motivating Example

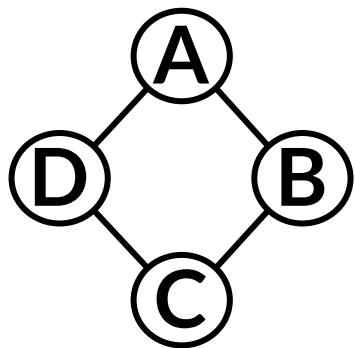
One way to define a probability over the joint voting decision of A, B, C, D is:

- Assign scores to each assignment of these variables, then define probability as the normalized score.
- A score can be any function, but we'll define it to be of the form:

$$\tilde{p}(A, B, C, D) = \phi(A, B) \phi(B, C) \phi(C, D) \phi(D, A)$$

- Where $\phi(X, Y)$ is a *factor* that assigns more weight to consistent vote among friends X, Y .
- For example:

$$\phi(X, Y) = \begin{cases} 10 & \text{if } X = Y = 1 \\ 5 & \text{if } X = Y = 0 \\ 1 & \text{otherwise.} \end{cases}$$

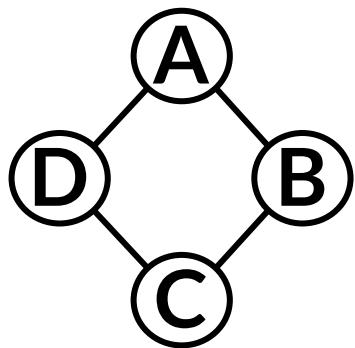


Markov Random Fields – Motivating Example

The final probability is then defined as

$$p(A, B, C, D) = \frac{1}{Z} \tilde{p}(A, B, C, D)$$

Where $Z = \sum_{A,B,C,D} \tilde{p}(A, B, C, D)$ is a *normalizing constant* that ensures the probability distribution sums to one.

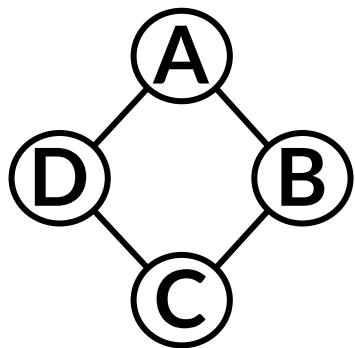


Markov Random Fields – Motivating Example

The final probability is then defined as

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

Where $Z = \sum_{A,B,C,D} \tilde{p}(A, B, C, D)$ is a *normalizing constant* that ensures the probability distribution sums to one.



Markov Random Fields – Motivating Example

The final probability is then defined as

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

Where $Z = \sum_{A,B,C,D} \tilde{p}(A, B, C, D)$ is a *normalizing constant* that ensures the probability distribution sums to one.

When normalized, we can view

- Factor $\phi(A, B)$ as an *interaction* pushing B's vote closer to that of A.
- And can view $\phi(B, C)$ as an *interaction* pushing B's vote closer to C, etc.
- And the most-likely vote will require reconciling these conflicting influences.

Markov Random Fields – Differences with Bayesian Networks

Note a few differences between Markov Random Fields and Bayesian Networks:

Markov Random Fields – Differences with Bayesian Networks

Note a few differences between Markov Random Fields and Bayesian Networks:

- In MRFs, we are not saying anything about how one variable is generated from another.
- Simply indicate a level of coupling between dependent variables in a graph.
- Don't have to specify a full generative story \Rightarrow could say it requires less prior knowledge.
- In turn this defines an *energy landscape* (unnormalized PDF), which is then converted into a probability (joint) distribution via the normalization constant.

Markov Random Fields – Formal Definition

Markov Random Fields – Formal Definition

Definition. A Markov Random Field (MRF) is a probability distribution p over variables x_1, \dots, x_n defined by an undirected graph G , with the form

$$p(x_1, \dots, x_n) = \frac{1}{Z} \prod_{c \in C} \phi_c(x_c)$$

where C denotes the set of cliques (i.e., fully connected subgraphs) of G , and each factor ϕ_c is a non-negative function over the variables in a clique.

Markov Random Fields – Formal Definition

Definition. A Markov Random Field (MRF) is a probability distribution p over variables x_1, \dots, x_n defined by an undirected graph G , with the form

$$p(x_1, \dots, x_n) = \frac{1}{Z} \prod_{c \in C} \phi_c(x_c)$$

where C denotes the set of cliques (i.e., fully connected subgraphs) of G , and each factor ϕ_c is a non-negative function over the variables in a clique.

And the *partition function* Z ensures that the distribution sums to one:

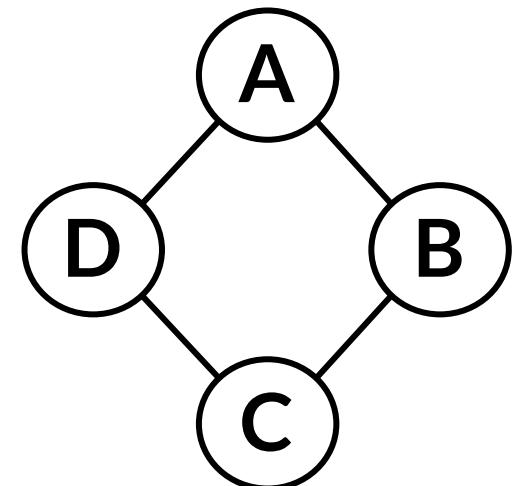
$$Z = \sum_{x_1, \dots, x_n} \prod_{c \in C} \phi_c(x_c)$$

Markov Random Fields – A few notes on cliques

Markov Random Fields – A few notes on cliques

Given a graph G , our probability distribution may contain factors whose scope is **any clique** in G .

Including: a single node, a pair of connected nodes, a triangle, etc.

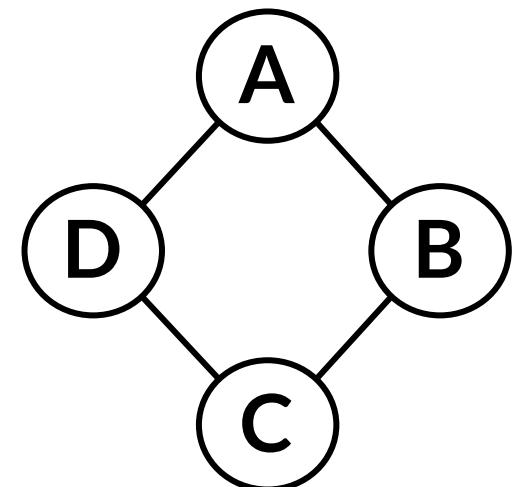


Markov Random Fields – A few notes on cliques

Given a graph G , our probability distribution may contain factors whose scope is **any clique** in G .

Including: a single node, a pair of connected nodes, a triangle, etc.

E.g., in friends example, could have specified factors defined on individual nodes (*unary factors*), but chose not to.

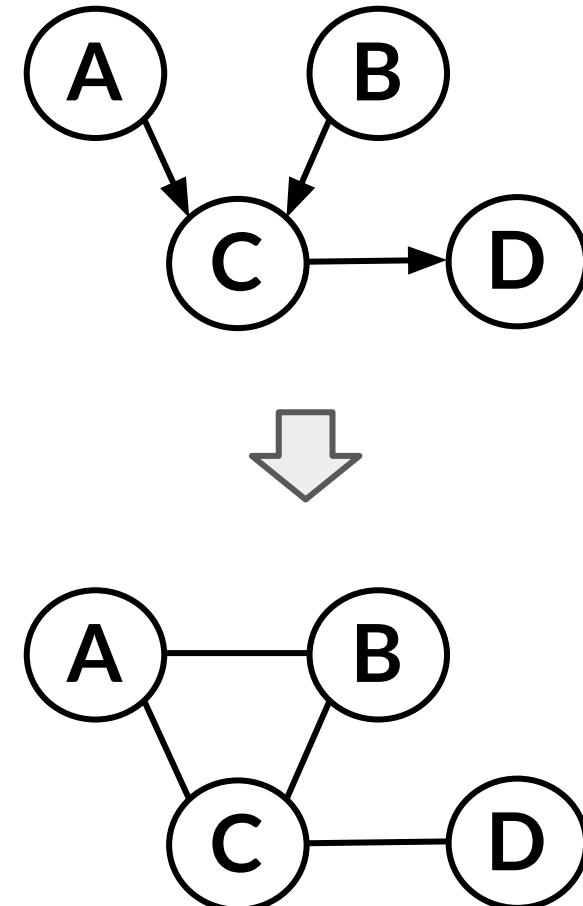


Markov Random Fields – Conversion Between BN and MRF

Markov Random Fields – Conversion Between BN and MRF

Bayesian networks are a special case of MRFs with a very specific type of clique factor.

⇒ Can convert from BN to MRF:

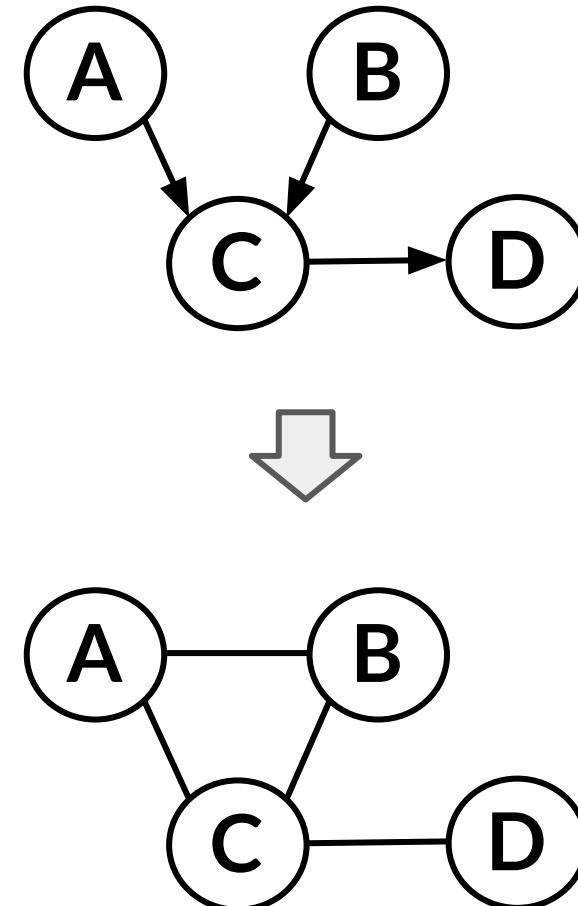


Markov Random Fields – Conversion Between BN and MRF

Bayesian networks are a special case of MRFs with a very specific type of clique factor.

⇒ Can convert from BN to MRF:

- If we take a directed graph G and add “side edges” to all parents of a given node (and removing their directionality)
- This is called “*moralization*”.

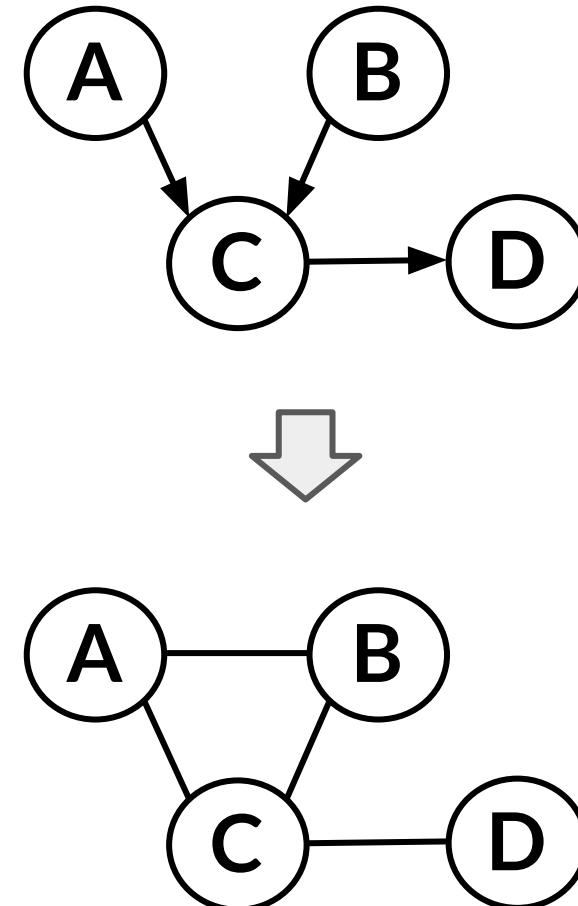


Markov Random Fields – Conversion Between BN and MRF

Bayesian networks are a special case of MRFs with a very specific type of clique factor.

⇒ Can convert from BN to MRF:

- If we take a directed graph G and add “side edges” to all parents of a given node (and removing their directionality)
- This is called “*moralization*”.
- Note: There is also a process for converting from MRF to BN, but you may lose some (conditional) independencies!



Markov Random Fields – Drawbacks

However, MRFs also have potential drawbacks:

Markov Random Fields – Drawbacks

However, MRFs also have potential drawbacks:

- Computing the normalization constant Z requires summing over a potentially exponential number of assignments.
- ⇒ In the general case, this is NP-hard; thus many undirected models will be intractable and will require approximation techniques.
- Undirected models may be difficult to interpret.
- It is much easier to generate data from a Bayesian network, which is important in some applications.

Some Famous Markov Random Fields

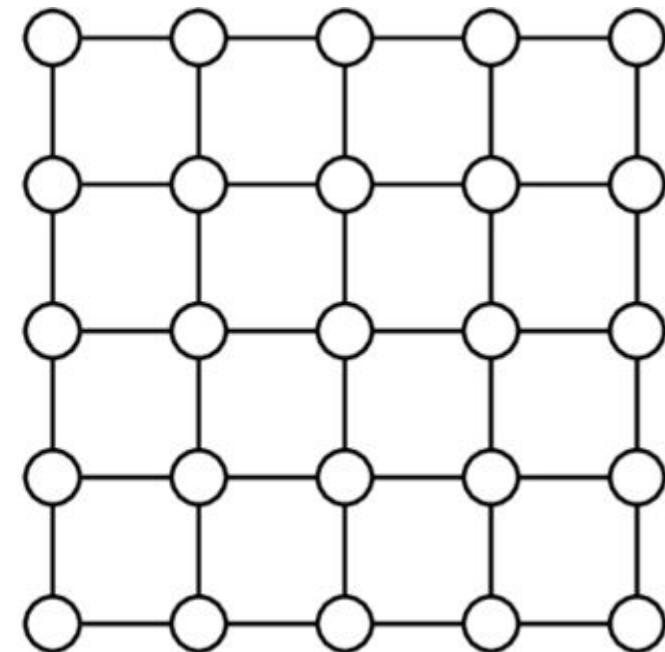
We'll go through a couple of famous MRFs

Some Famous Markov Random Fields – Ising Model

Some Famous Markov Random Fields – Ising Model

Originally a model of ferromagnetism in statistical mechanics!

Markov Random Field



Some Famous Markov Random Fields – Ising Model

Originally a model of ferromagnetism in statistical mechanics!

MRF is very simple:

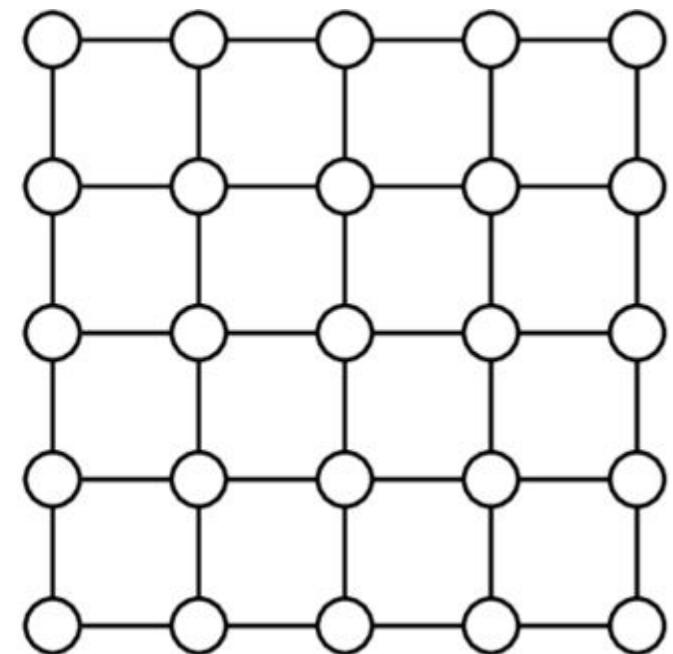
A lattice of nodes (“*grid graph*”).

Each typically a binary random variable.

- Represent “spins” that can be in one of two states.

Neighboring spins that agree have a lower energy than those that disagree.

Markov Random Field

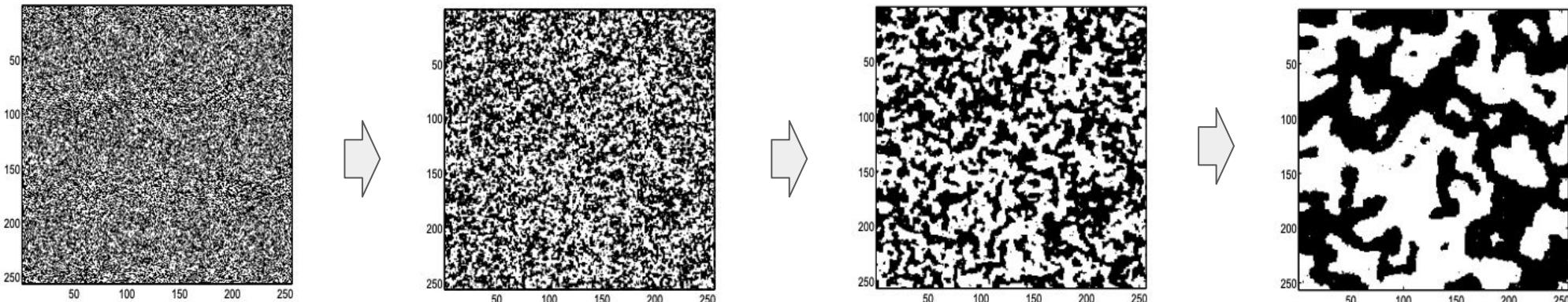


Some Famous Markov Random Fields – Ising Model

Simulating samples from the joint PDF.

Some Famous Markov Random Fields – Ising Model

Simulating samples from the joint PDF.



Source: Brani Vidakovic, "Markov Random Fields"

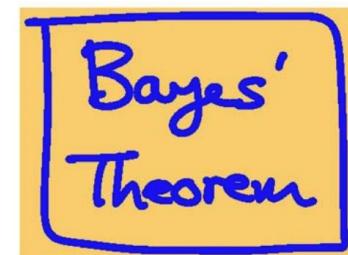
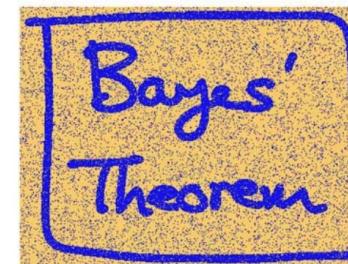
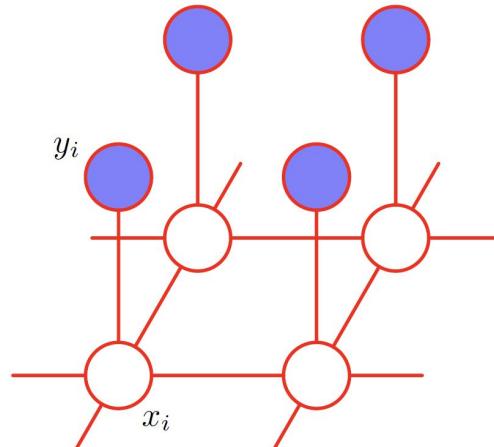
⇒ Neighbors tend to share values!

Some Famous Markov Random Fields – Ising Model

Applications:

Some Famous Markov Random Fields – Ising Model

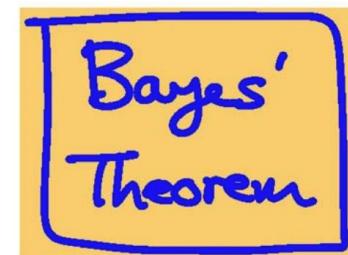
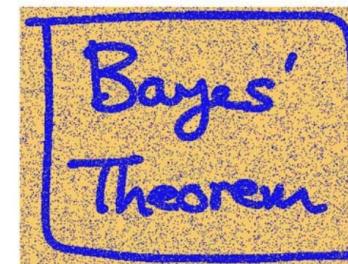
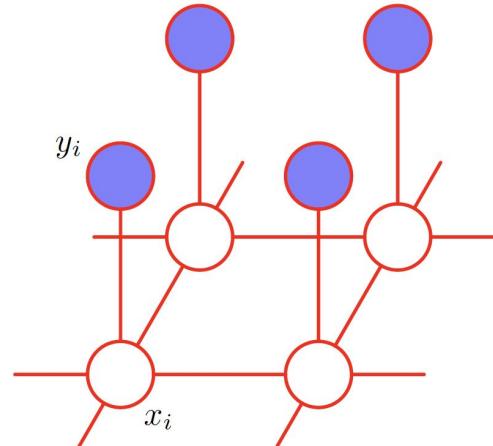
Applications: image denoising.



Source: Sargur Srihari, "Markov Random Fields"

Some Famous Markov Random Fields – Ising Model

Applications: image denoising.



Source: Sargur Srihari, "Markov Random Fields"

Along with models of: phase separation, liquid-gas transitions, binary alloys, spin glasses, protein folding, opinion dynamics, social network analysis, etc.

Some Famous Markov Random Fields – Conditional Random Field

Some Famous Markov Random Fields – Conditional Random Field

Conditional random field – similar to an undirected version of a HMM.

- Predict latent variables based on observed variables (on sequences).

Original Paper (2001)

→ 19,154 citations

Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data

John Lafferty^{1*}
Andrew McCallum^{1†}
Fernando Pereira^{1‡}
¹WhizBang! Labs-Research, 4616 Henry Street, Pittsburgh, PA 15213 USA
¹School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213 USA
¹Department of Computer and Information Science, University of Pennsylvania, Philadelphia, PA 19104 USA

LAFFERTY@CS.CMU.EDU
MCCALLUM@WHIZBANG.COM
FPEREIRA@WHIZBANG.COM

Abstract

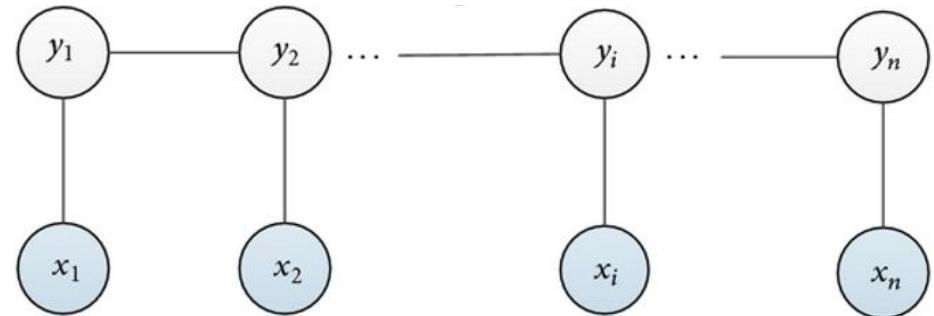
We present *conditional random fields*, a framework for building probabilistic models to segment and label sequence data. Conditional random fields offer several advantages over hidden Markov models and stochastic grammars for these tasks, including the ability to relax strong independence assumptions made in those models. Conditional random fields also avoid a fundamental limitation of maximum entropy Markov models (MEMMs) and other discriminative Markov models based on directed graphical models, which can be biased towards states with few successor states. We present iterative parameter estimation algorithms for conditional random fields and compare the performance of the resulting models to HMMs and MEMMs on synthetic and natural-language data.

1. Introduction

The need to segment and label sequences arises in many different problems in several scientific fields. Hidden Markov models (HMMs) are the most well understood and widely used probabilistic models for such problems. In computational biology, HMMs and stochastic grammars have been successfully used to align biological sequences, find sequences homologous to a known evolutionary family, and analyze RNA secondary structure (Durbin et al., 1998). In computational linguistics and computer science, HMMs and stochastic grammars have been applied to a wide variety of problems in text and speech processing, including topic segmentation, part-of-speech (POS) tagging, information extraction, and syntactic disambiguation (Manning & Schütze, 1999).

Maximum entropy Markov models (MEMMs) are conditionally probabilistic sequence models that attain all of the above advantages (McCallum et al., 2000). In MEMMs, each state “state” has an exponential model that takes the observation sequence as input, and outputs a distribution over possible next states. These exponential models are trained by an appropriate iterative scaling method in the

*Output labels are associated with states; it is possible for several states to have the same label. To, for simplicity, in the rest of this paper we assume a one-to-one correspondence.



Source: "Conditional Random Fields", codersarts.com

Some Famous Markov Random Fields – Conditional Random Field

Conditional random field – similar to an undirected version of a HMM.

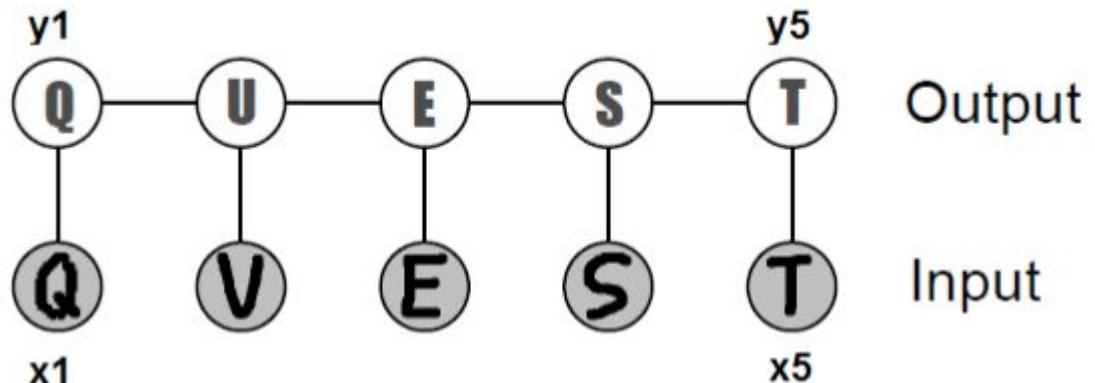
- Predict latent variables based on observed variables (on sequences).

Original Paper (2001)

→ 19,154 citations



E.g., used to segment or label sequence data.



Source: Stefano Ermon, Deep Generative Models (CS236) Class

Applications in NLP, computer vision, bioinformatics.

Conclusion

Conclusion

What are PGMS? A rich visual language for defining conditional independence assumptions (and other structure) in joint probability distributions.

Conclusion

What are PGMS? A rich visual language for defining conditional independence assumptions (and other structure) in joint probability distributions.

Today:

- Directed vs undirected graphical models: Bayesian networks and Markov random fields.
- Plate notation, generative process notation.
- Observed variables vs. latent variables vs. (hyper)parameters.
- Famous/classic Bayes Nets (HMMs, GMMs, LDA, VAE, etc.)
- Pros and Cons of Bayes Nets vs MRFs.
- Famous/classic MRFs (Ising Model, CRF).

Assignments and Grading – Grading Breakdown

| <u>Assignment</u> | <u>% of Grade</u> |
|--|-------------------|
| 1. Paper Presentation | 20% |
| 2. In-class Participation and Discussion | |
| 2a. Role 1 – Discussion Lead 1 | 8% |
| 2a. Role 2 – Discussion Lead 2 | 8% |
| 2a. Role 3 – Scribe | 9% |
| 3. Course Project | |
| 3a. Project Pitch | 8% |
| 3b. Midway Report | 10% |
| 3c. Final Presentation | 12% |
| 3d. Final Report | 25% |

Course Project – Group Project

This will be a **group project** – groups of 3-4 students.

- Aiming for ~10 groups total (due to timing constraints)
- We will help facilitate this during class.
- *E.g.*, everyone will introduce themselves and describe research interests, which we will write/share, to help in matching.
- Need to aim to form teams and select project idea by roughly end of this month.

Course Project – Guidance & Expectations

What does this project entail?

Course Project – Guidance & Expectations

What does this project entail?

I want people to use a probabilistic or generative model in some way!

- Application of prob/gen models from this class on a novel task or dataset.
- Algorithmic improvements in learning, inference, or evaluation of prob/gen models.
- Theoretical analysis of any aspect of existing prob/gen models.

Course Project – Guidance & Expectations

What does this project entail?

I want people to use a probabilistic or generative model in some way!

- Application of prob/gen models from this class on a novel task or dataset.
- Algorithmic improvements in learning, inference, or evaluation of prob/gen models.
- Theoretical analysis of any aspect of existing prob/gen models.

Goal is to complete a small-scale implementation or pilot study during the class.
I'm more focused on interesting conceptual ideas, rather than on performance/results.

Aim to connect it to the research you are focusing on outside of this class!

Course Project – Assignments

The project will be worth a substantial portion of the grade, and consist of four main assignments:

Course Project – Assignments

The project will be worth a substantial portion of the grade, and consist of four main assignments:

- **Project pitch:** Each group will come up with a project idea, make a few slides, and share their idea with the class for feedback (**10 minutes long, on Feb 7 & Feb 14**)

Course Project – Assignments

The project will be worth a substantial portion of the grade, and consist of four main assignments:

- **Project pitch:** Each group will come up with a project idea, make a few slides, and share their idea with the class for feedback (**10 minutes long, on Feb 7 & Feb 14**)
- **Midway report:** Each group will write a short (4 page) midway report for their project, focusing on a literature review, implementation plan, and any initial experiments.
 - Latex template will be provided.

Course Project – Assignments

The project will be worth a substantial portion of the grade, and consist of four main assignments:

- **Project pitch:** Each group will come up with a project idea, make a few slides, and share their idea with the class for feedback (**10 minutes long, on Feb 7 & Feb 14**)
- **Midway report:** Each group will write a short (4 page) midway report for their project, focusing on a literature review, implementation plan, and any initial experiments.
 - Latex template will be provided.
- **Final presentation:** Each group will give a presentation to the class on their final project (**30 minutes long, on Apr 25 & May 2**)

Course Project – Assignments

The project will be worth a substantial portion of the grade, and consist of four main assignments:

- **Project pitch:** Each group will come up with a project idea, make a few slides, and share their idea with the class for feedback (**10 minutes long, on Feb 7 & Feb 14**)
- **Midway report:** Each group will write a short (4 page) midway report for their project, focusing on a literature review, implementation plan, and any initial experiments.
 - Latex template will be provided.
- **Final presentation:** Each group will give a presentation to the class on their final project (**30 minutes long, on Apr 25 & May 2**)
- **Final report:** Each group will submit a final report for their project, describing all details, background, prior work, and results (**8-10 pages long, due May 9**).
 - Latex template will be provided.

Student Introductions

Link to spreadsheet:

<https://docs.google.com/spreadsheets/d/1h5VtS7vFTn8mgsYYBQNDIpBxFI6OUKbrUKNrtXF6iDg/edit?usp=sharing>