

---

---

# Advanced Prediction Models

Deep Learning, Graphical Models and Reinforcement  
Learning

# Recap: Why Graphical Models

---

- We have seen deep learning techniques for unstructured data
  - Predominantly vision and text/audio
  - We will see control in the last part of the course
    - (Reinforcement Learning)
- For structured data, graphical models are the most versatile framework
  - Successfully applications:
    - Kalman filtering in engineering
    - Decoding in cell phones (channel codes)
    - Hidden Markov models for time series
    - Clustering, regression, classification ...

# Recap: Graphical Models Landscape

---

- Three key parts:
  - Representation
    - Capture uncertainty (joint distribution)
    - Capture **conditional independences** (metadata)
    - Visualization of metadata for a distribution
  - Inference
    - Efficient methods for computing marginal or conditional distributions **quickly**
  - Learning
    - Learning the **parameters of the distribution** can deal with prior knowledge and missing data

# Today's Outline

---

- Applications
- Learning
  - DPGM/UPGM
  - Parameter Estimation
  - Structure Estimation
  - Complete/Incomplete Data

---

---

# Applications

# Applications of Graphical Models

---

- Given all that we have learned up to now, we will sample the following applications

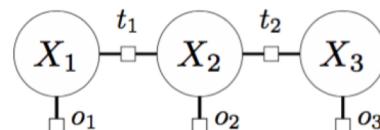
<b>Hidden Markov Models</b>	time series, tracking
<b>Gaussian Mixture Models</b>	clustering
<b>Latent Dirichlet Allocation</b>	topic modeling
<b>Conditional Random Fields</b>	structured classification/regression

# Example Graphical Model I



## Problem: person tracking

Sensors reports positions: 0, 2, 2. Objects don't move very fast and sensors are a bit noisy. What path did the person take?



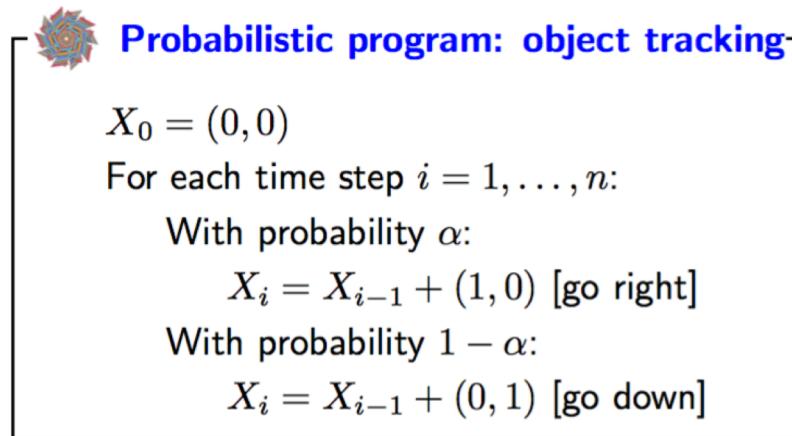
- Variables  $X_i$ : location of object at position  $i$
- Transition factors  $t_i(x_i, x_{i+1})$ : incorporate physics
- Observation factors  $o_i(x_i)$ : incorporate sensors

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

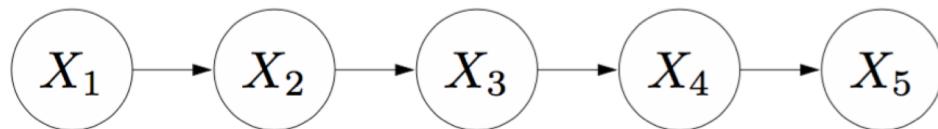
# Example Graphical Model II

- A generative process is nothing but a description of the joint distribution in terms of how the random variables realize

Probabilistic program:



Bayesian network:



Mathematical definition:

$$p(x_i | x_{i-1}) = \alpha \cdot \underbrace{[x_i = x_{i-1} + (1, 0)]}_{\text{right}} + (1 - \alpha) \cdot \underbrace{[x_i = x_{i-1} + (0, 1)]}_{\text{down}}$$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Example Graphical Model III

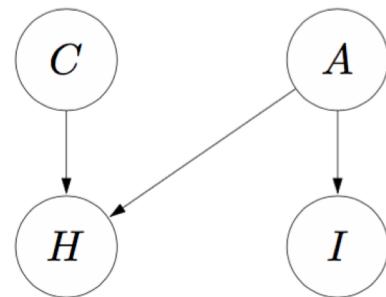


Problem: cold or allergies?

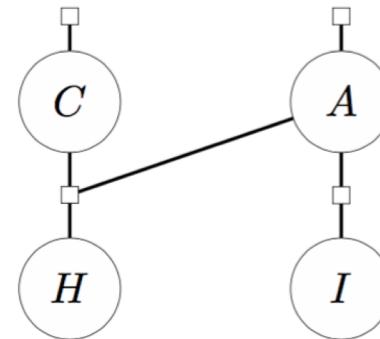
You are coughing and have itchy eyes. What do you have?

Variables: Cold, Allergy, Cough, Itchy eyes

Bayesian network:



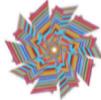
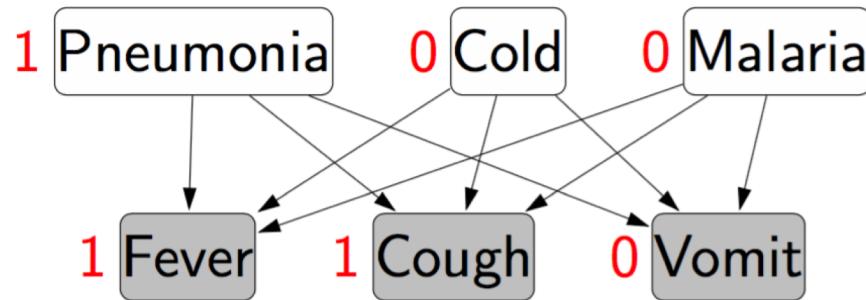
Factor graph:



<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Example Graphical Model IV

Question: If patient has has a cough and fever, what disease(s) does he/she have?



## Probabilistic program: diseases and symptoms

For each disease  $i = 1, \dots, m$ :

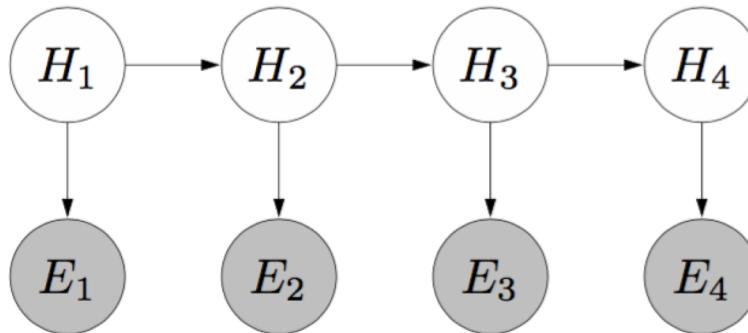
Generate activity of disease  $D_i \sim p(D_i)$

For each symptom  $j = 1, \dots, n$ :

Generate activity of symptom  $S_j \sim p(S_j | D_{1:m})$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Object Tracking via Hidden Markov Model



## Problem: object tracking

$H_i \in \{1, \dots, K\}$ : location of object at time step  $i$

$E_i \in \{1, \dots, K\}$ : sensor reading at time step  $i$

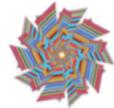
Start:  $p(h_1)$ : uniform over all locations

Transition  $p(h_i | h_{i-1})$ : uniform over adjacent loc.

Emission  $p(e_i | h_i)$ : uniform over adjacent loc.

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Generative Program for HMM

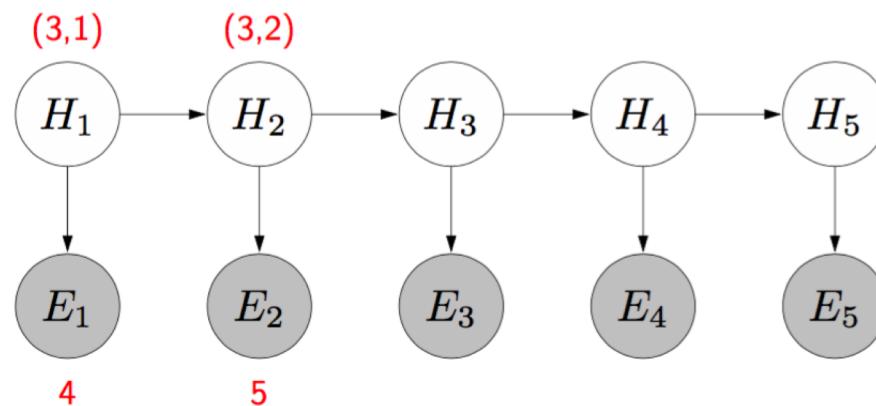


## Probabilistic program: hidden Markov model (HMM)

For each time step  $t = 1, \dots, T$ :

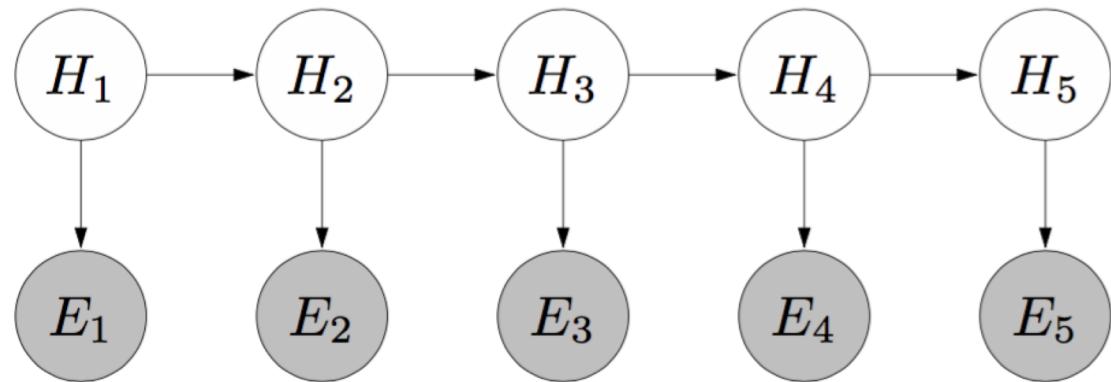
Generate object location  $H_t \sim p(H_t | H_{t-1})$

Generate sensor reading  $E_t \sim p(E_t | H_t)$



<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Object Tracking via HMM



$$\mathbb{P}(H = h, E = e) = \underbrace{p(h_1)}_{\text{start}} \prod_{i=2}^n \underbrace{p(h_i | h_{i-1})}_{\text{transition}} \prod_{i=1}^n \underbrace{p(e_i | h_i)}_{\text{emission}}$$

Query (**filtering**):

$$\mathbb{P}(H_3 | E_1 = e_1, E_2 = e_2, E_3 = e_3)$$

Query (**smoothing**):

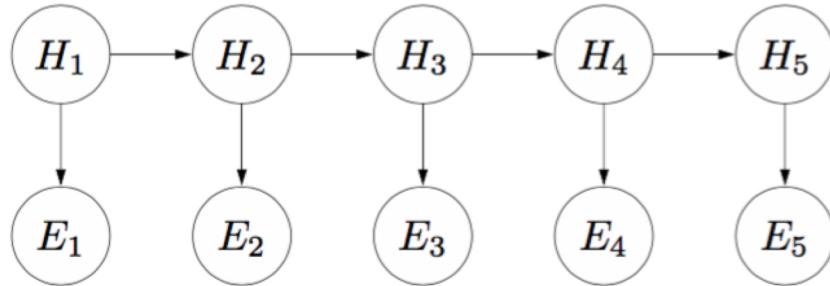
$$\mathbb{P}(H_3 | E_1 = e_1, E_2 = e_2, E_3 = e_3, E_4 = e_4, E_5 = e_5)$$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# HMM Parameter Sharing

Variables:

- $H_1, \dots, H_n$  (e.g., actual positions)
- $E_1, \dots, E_n$  (e.g., sensor readings)



$$\mathbb{P}(H = h, E = e) = \prod_{i=1}^n p_{\text{trans}}(h_i | h_{i-1}) p_{\text{emit}}(e_i | h_i)$$

Parameters:  $\theta = (p_{\text{trans}}, p_{\text{emit}})$

$\mathcal{D}_{\text{train}}$  is a set of full assignments to  $(H, E)$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Mixture Models

---

- “Standard” distributions (e.g., multivariate Gaussian) are too limited
- How do we represent and learn more complex ones?
- One answer: Mixtures of “standard” distributions
- In the limit, can approximate any distribution this way
- Also good (and widely used) as a clustering method

---

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Gaussian Mixture Models

---

- The  $N$ -dim. multivariate normal distribution,  $\mathcal{N}(\mu, \Sigma)$ , has density:

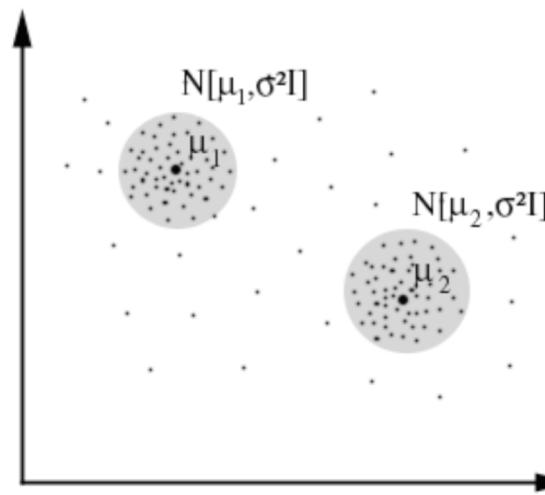
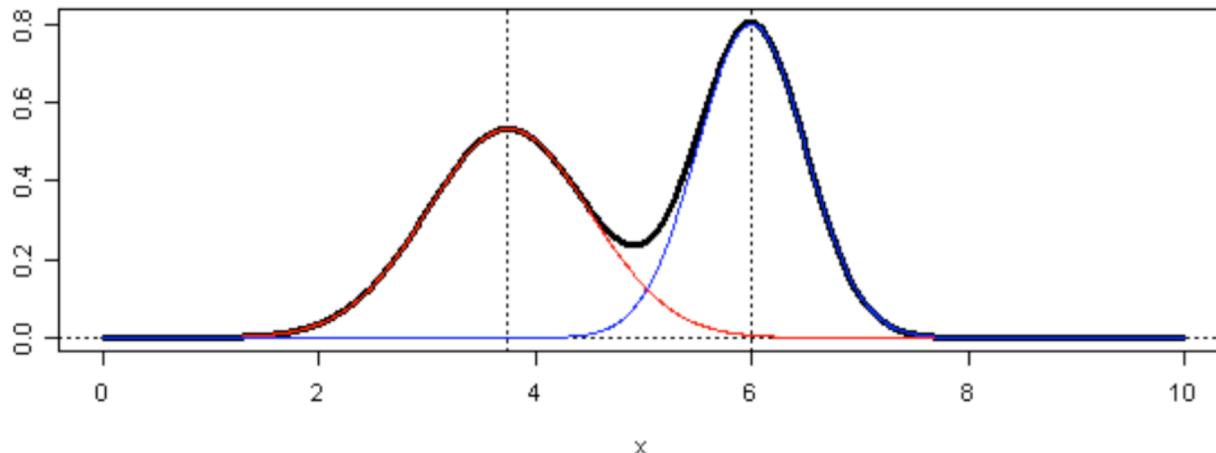
$$p(\mathbf{x}) = \frac{1}{(2\pi)^{N/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right)$$

- Suppose we have  $k$  Gaussians given by  $\mu_k$  and  $\Sigma_k$ , and a distribution  $\theta$  over the numbers  $1, \dots, k$
- Mixture of Gaussians distribution  $p(y, \mathbf{x})$  given by
  - ① Sample  $y \sim \theta$  (specifies which Gaussian to use)
  - ② Sample  $\mathbf{x} \sim \mathcal{N}(\mu_y, \Sigma_y)$

<sup>1</sup>Reference: David Sontag (2013)

# Gaussian Mixture Model in 1D and 2D

- The marginal distribution over  $x$  looks like:



# Learning a GMM

---

Initialize parameters ignoring missing information

Repeat until convergence:

**E step:** Compute expected values of unobserved variables, assuming current parameter values

**M step:** Compute new parameter values to maximize probability of data (observed & estimated)

(Also: Initialize expected values ignoring missing info)

# Learning a 1D GMM

**Initialization:** Choose means at random, etc.

**E step:** For all examples  $x_k$ :

$$P(\mu_i|x_k) = \frac{P(\mu_i)P(x_k|\mu_i)}{P(x_k)} = \frac{P(\mu_i)P(x_k|\mu_i)}{\sum_{i'} P(\mu_{i'})P(x_k|\mu_{i'})}$$

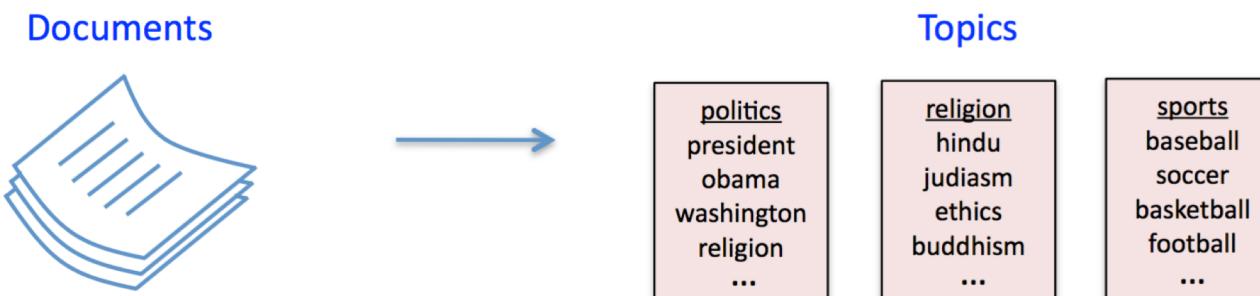
**M step:** For all components  $c_i$ :

$$\begin{aligned} P(c_i) &= \frac{1}{n_e} \sum_{k=1}^{n_e} P(\mu_i|x_k) \\ \mu_i &= \frac{\sum_{k=1}^{n_e} x_k P(\mu_i|x_k)}{\sum_{k=1}^{n_e} P(\mu_i|x_k)} \\ \sigma_i^2 &= \frac{\sum_{k=1}^{n_e} (x_k - \mu_i)^2 P(\mu_i|x_k)}{\sum_{k=1}^{n_e} P(\mu_i|x_k)} \end{aligned}$$

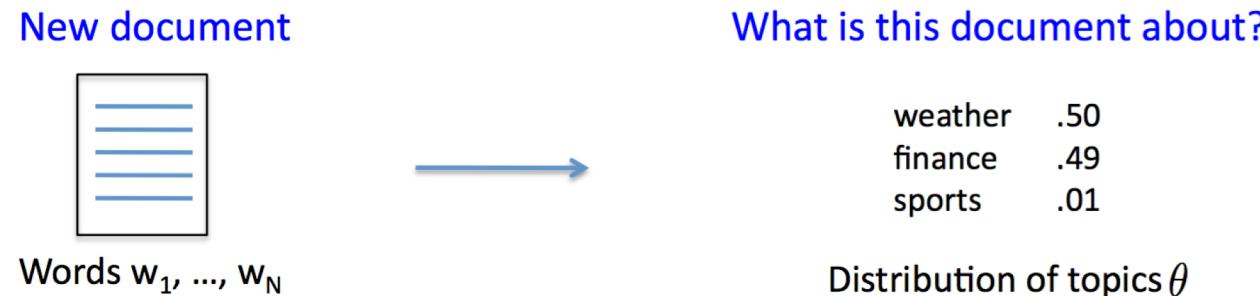
<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Latent Dirichlet Allocation

- **Topic models** are powerful tools for exploring large data sets and for making inferences about the content of documents



- Many applications in information retrieval, document summarization, and classification



- LDA is one of the simplest and most widely used topic models

<sup>1</sup>Reference: David Sontag (2013)

# Latent Dirichlet Allocation

---

- ➊ Sample the document's **topic distribution**  $\theta$  (aka topic vector)

$$\theta \sim \text{Dirichlet}(\alpha_{1:T})$$

where the  $\{\alpha_t\}_{t=1}^T$  are fixed hyperparameters. Thus  $\theta$  is a distribution over  $T$  topics with mean  $\theta_t = \alpha_t / \sum_t \alpha_t$

- ➋ For  $i = 1$  to  $N$ , sample the **topic**  $z_i$  of the  $i$ 'th word

$$z_i | \theta \sim \theta$$

- ➌ ... and then sample the actual **word**  $w_i$  from the  $z_i$ 'th topic

$$w_i | z_i \sim \beta_{z_i}$$

where  $\{\beta_t\}_{t=1}^T$  are the *topics* (a fixed collection of distributions on words)

<sup>1</sup>Reference: David Sontag (2013)

# Latent Dirichlet Allocation

... and then sample the actual **word**  $w_i$  from the  $z_i$ 'th topic

$$w_i | z_i \sim \beta_{z_i}$$

where  $\{\beta_t\}_{t=1}^T$  are the *topics* (a fixed collection of distributions on words)

Documents



Topics

<u>politics</u> .0100
president .0095
obama .0090
washington .0085
religion .0060
...

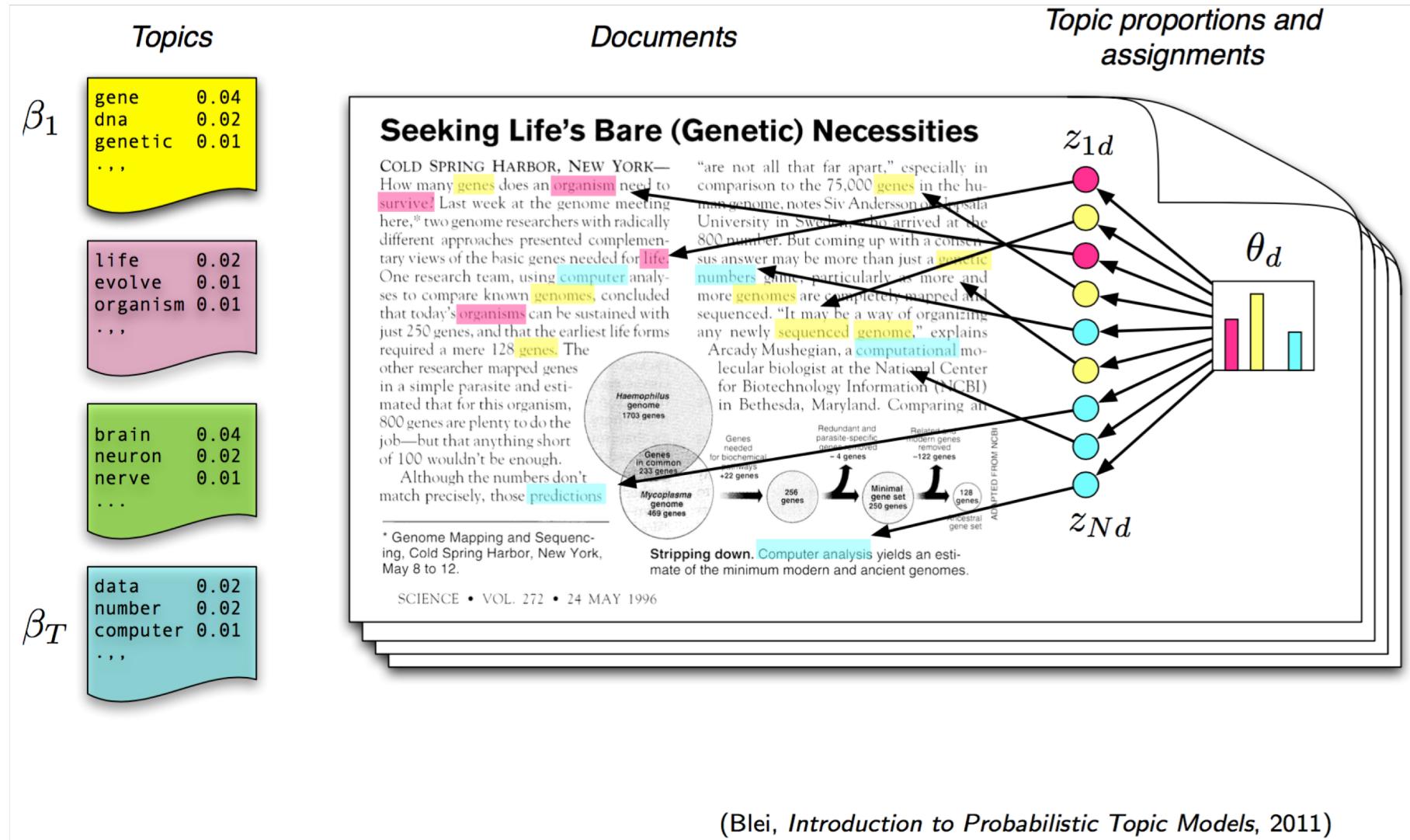
<u>religion</u> .0500
hindu .0092
judiasm .0080
ethics .0075
buddhism .0016
...

<u>sports</u> .0105
baseball .0100
soccer .0055
basketball .0050
football .0045
...

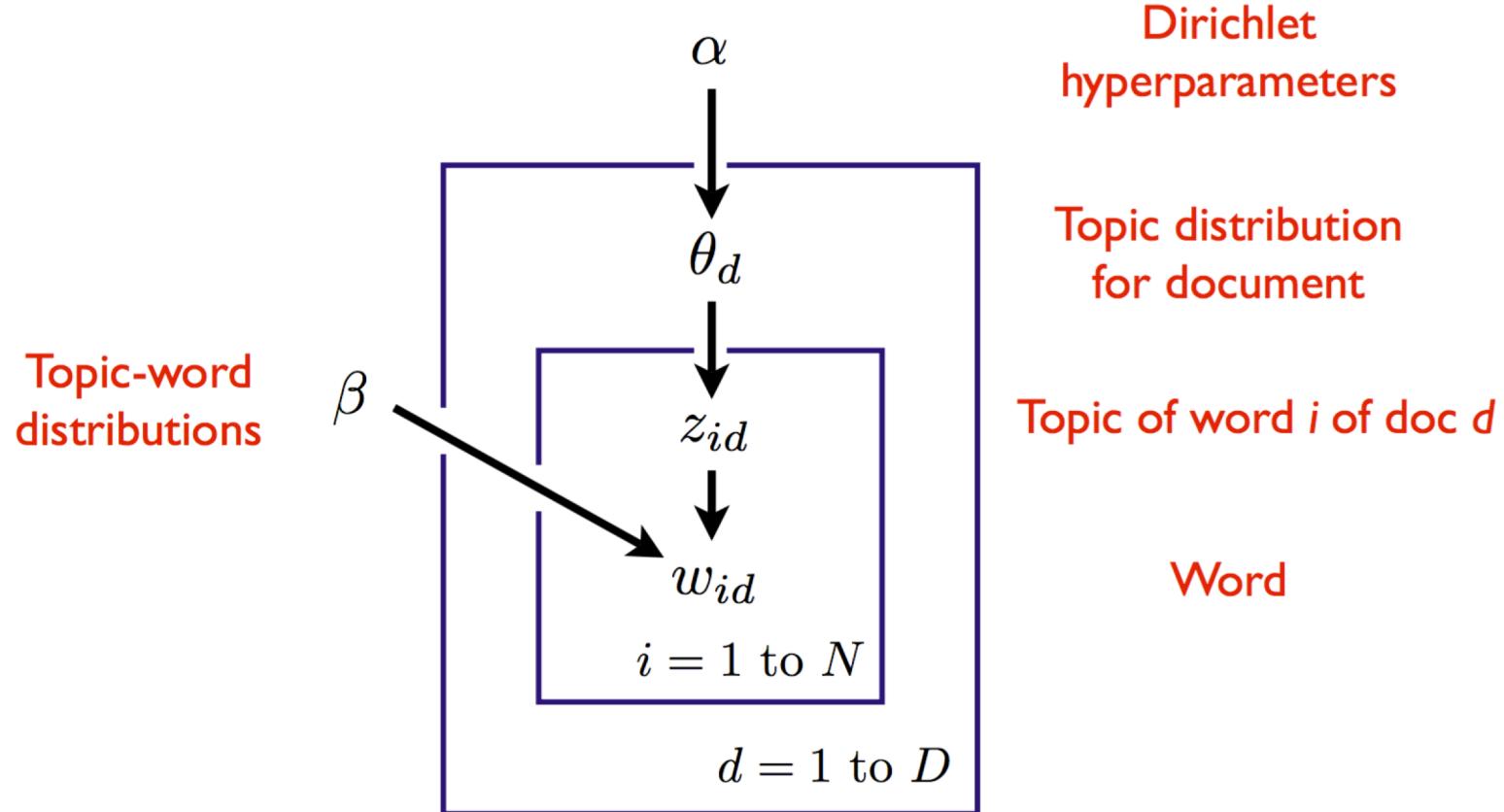
$$\beta_t = \{ p(w | z = t) \}$$

<sup>1</sup>Reference: David Sontag (2013)

# Latent Dirichlet Allocation



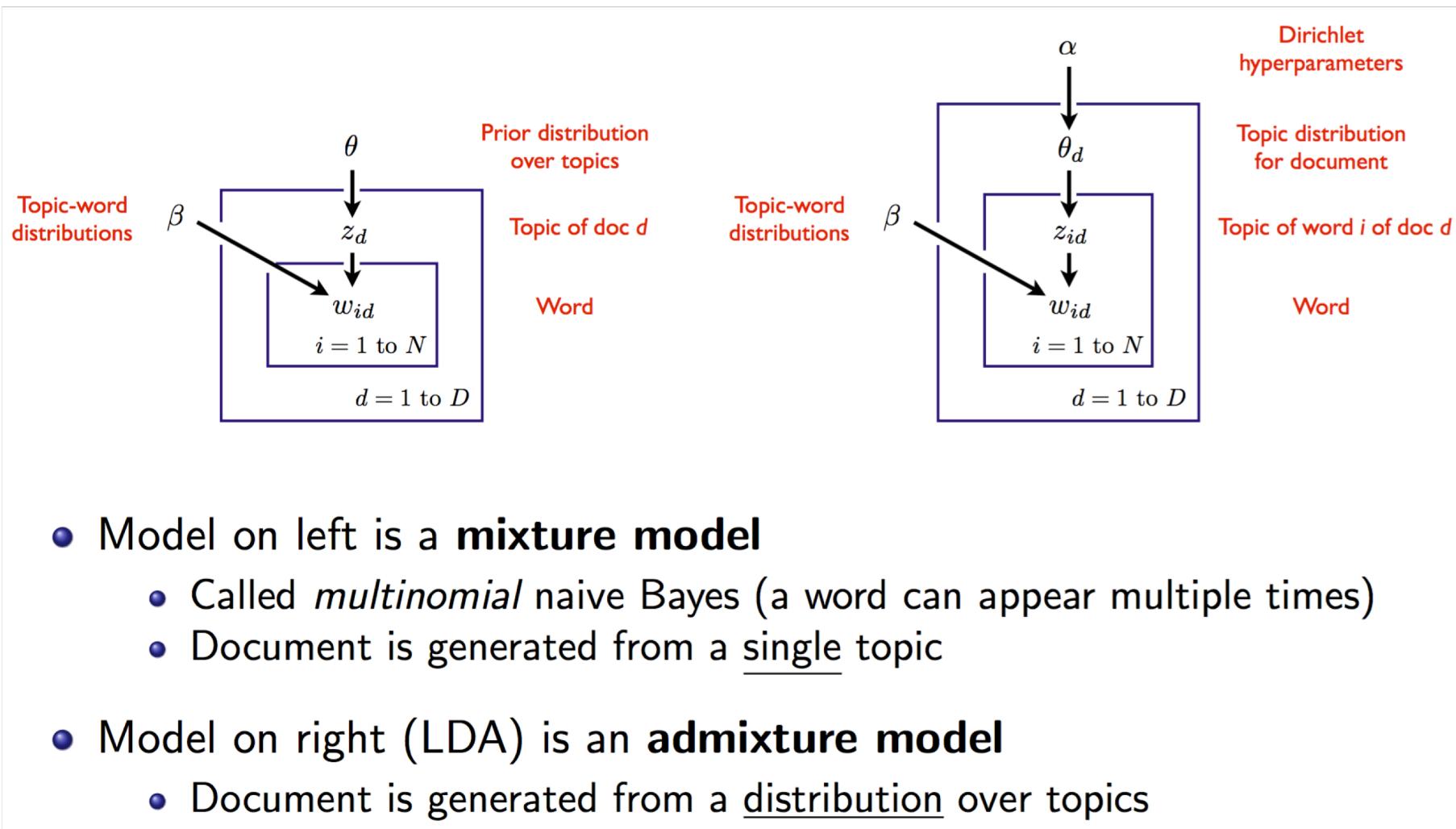
# Latent Dirichlet Allocation



Variables within a plate are replicated in a conditionally independent manner

<sup>1</sup>Reference: David Sontag (2013)

# Latent Dirichlet Allocation



- Model on left is a **mixture model**
  - Called *multinomial* naive Bayes (a word can appear multiple times)
  - Document is generated from a single topic
- Model on right (LDA) is an **admixture model**
  - Document is generated from a distribution over topics

<sup>1</sup>Reference: David Sontag (2013)

# Conditional Random Field based Classifier

---

- **Conditional random fields** are undirected graphical models of conditional distributions  $p(\mathbf{Y} | \mathbf{X})$ 
  - $\mathbf{Y}$  is a set of **target variables**
  - $\mathbf{X}$  is a set of **observed variables**
- We typically show the graphical model using just the  $\mathbf{Y}$  variables
- Potentials are a function of  $\mathbf{X}$  and  $\mathbf{Y}$

<sup>1</sup>Reference: David Sontag (2013)

# Conditional Random Field based Classifier

---

- A CRF is a Markov network on variables  $\mathbf{X} \cup \mathbf{Y}$ , which specifies the conditional distribution

$$P(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{c \in C} \phi_c(\mathbf{x}_c, \mathbf{y}_c)$$

with partition function

$$Z(\mathbf{x}) = \sum_{\hat{\mathbf{y}}} \prod_{c \in C} \phi_c(\mathbf{x}_c, \hat{\mathbf{y}}_c).$$

- As before, two variables in the graph are connected with an undirected edge if they appear together in the scope of some factor
- The only difference with a standard Markov network is the normalization term – before marginalized over  $\mathbf{X}$  and  $\mathbf{Y}$ , now only over  $\mathbf{Y}$

<sup>1</sup>Reference: David Sontag (2013)

# CRF for Natural Language Processing: Log-linear Terms

---

- Factors may depend on a large number of variables
- We typically parameterize each factor as a log-linear function,

$$\phi_c(\mathbf{x}_c, \mathbf{y}_c) = \exp\{\mathbf{w} \cdot \mathbf{f}_c(\mathbf{x}_c, \mathbf{y}_c)\}$$

- $\mathbf{f}_c(\mathbf{x}_c, \mathbf{y}_c)$  is a feature vector
- $\mathbf{w}$  is a weight vector which is typically learned – we will discuss this extensively in later lectures

<sup>1</sup>Reference: David Sontag (2013)

# CRF for Natural Language Processing: The Task

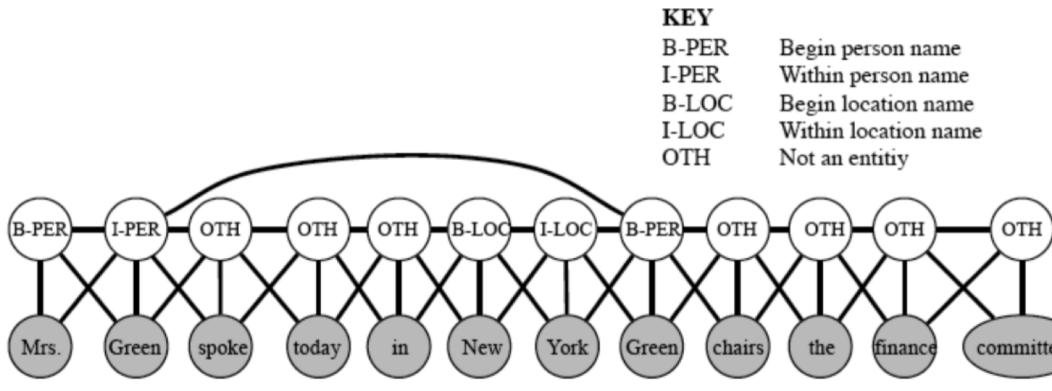
---

- Given a sentence, determine the people and organizations involved and the relevant locations:  
“Mrs. Green spoke today in New York. Green chairs the finance committee.”
- Entities sometimes span multiple words. Entity of a word not obvious without considering its *context*
- CRF has one variable  $X_i$  for each word, and  $Y_i$  encodes the possible labels of that word
- The labels are, for example, “B-person, I-person, B-location, I-location, B-organization, I-organization”
  - Having beginning (B) and within (I) allows the model to segment adjacent entities

<sup>1</sup>Reference: David Sontag (2013)

# CRF for Natural Language Processing: The Task

The graphical model looks like (called a *skip-chain CRF*):



There are three types of potentials:

- $\phi^1(Y_t, Y_{t+1})$  represents dependencies between neighboring target variables [analogous to transition distribution in a HMM]
- $\phi^2(Y_t, Y_{t'})$  for all pairs  $t, t'$  such that  $x_t = x_{t'}$ , because if a word appears twice, it is likely to be the same entity
- $\phi^3(Y_t, X_1, \dots, X_T)$  for dependencies between an entity and the word sequence [e.g., may have features taking into consideration capitalization]

**Notice that the graph structure changes depending on the sentence!**

<sup>1</sup>Reference: David Sontag (2013)

---

---

# Questions?

# Today's Outline

---

- Applications
- Learning
  - Parameter Estimation in DPGMs with Complete/Incomplete Data
  - Structure Estimation in DPGMs
  - Parameter Estimation in UPGMs with Complete/Incomplete Data

---

---

# Estimation/Learning

# Different Estimation/Learning Problems

---

- There are many variants

<b>Model</b>	DPGM	UPGM
<b>Data</b>	Complete	Incomplete
<b>Structure</b>	Known	Unknown
<b>Objective</b>	Generative	Discriminative

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Different Estimation/Learning Problems

---

- We will look at the following problems
  - Learning DPGMs with complete data and known structure
    - MLE via counting and normalizing
  - Learning DPGMs with incomplete data and known structure
    - EM
  - Learning DPGM structure
  - Learning UPGMs in a generative setting
  - Learning UPGM in a discriminative setting

# Different Estimation/Learning Problems

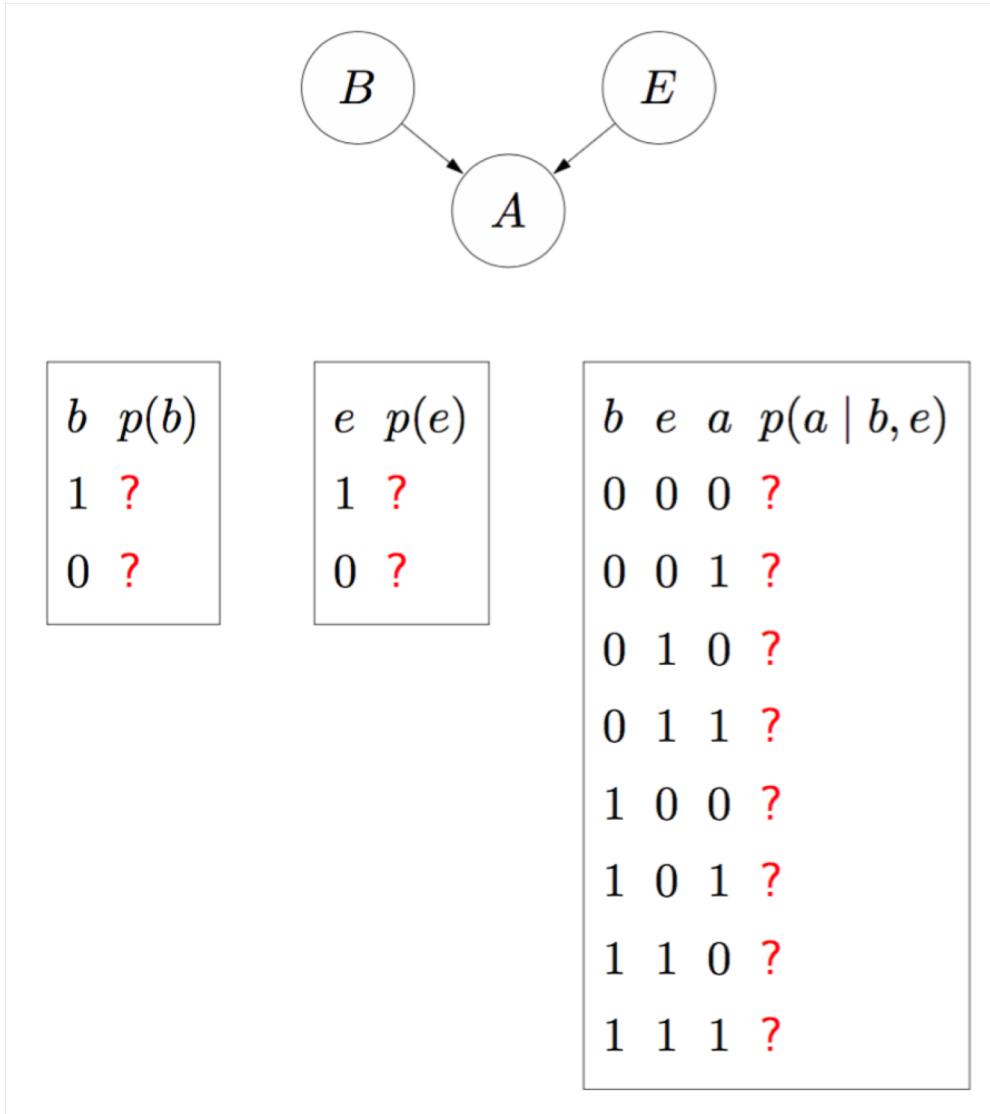
---

- There are many variants

<b>Model</b>	DPGM	UPGM
<b>Data</b>	Complete	Incomplete
<b>Structure</b>	Known	Unknown
<b>Objective</b>	Generative	Discriminative

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Learning in DPGM: Parameters



<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: Parameter Estimation

---

**Training data**

$\mathcal{D}_{\text{train}}$  (an example is an assignment to  $X$ )



**Parameters**

$\theta$  (local conditional probabilities)

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: One Variable Example

---

Setup:

- One variable  $R$  representing the rating of a movie  $\{1, 2, 3, 4, 5\}$

$$R \quad \mathbb{P}(R = r) = p(r)$$

Parameters:

$$\theta = (p(1), p(2), p(3), p(4), p(5))$$

Training data:

$$\mathcal{D}_{\text{train}} = \{1, 3, 4, 4, 4, 4, 4, 4, 5, 5, 5\}$$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: One Variable Example

Learning:

$$\mathcal{D}_{\text{train}} \Rightarrow \theta$$

Intuition:  $p(r) \propto$  number of occurrences of  $r$  in  $\mathcal{D}_{\text{train}}$

Example:

$$\mathcal{D}_{\text{train}} = \{1, 3, 4, 4, 4, 4, 4, 5, 5, 5\}$$



$\theta$ :

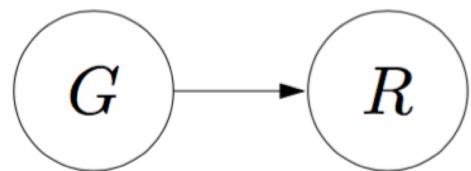
$r$	$p(r)$
1	0.1
2	0.0
3	0.1
4	0.5
5	0.3

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: Two Variables Example

Variables:

- Genre  $G \in \{\text{drama, comedy}\}$
- Rating  $R \in \{1, 2, 3, 4, 5\}$



$$\mathbb{P}(G = g, R = r) = p_G(g)p_R(r | g)$$

$$\mathcal{D}_{\text{train}} = \{(d, 4), (d, 4), (d, 5), (c, 1), (c, 5)\}$$

Parameters:  $\theta = (p_G, p_R)$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: Two Variables Example

$$\mathcal{D}_{\text{train}} = \{(d, 4), (d, 4), (d, 5), (c, 1), (c, 5)\}$$

Intuitive strategy:

- Estimate each local conditional distribution ( $p_G$  and  $p_R$ ) separately
- For each value of conditioned variable (e.g.,  $g$ ), estimate distribution over values of unconditioned variable (e.g.,  $r$ )

$\theta:$

$g$	$p_G(g)$
d	3/5
c	2/5

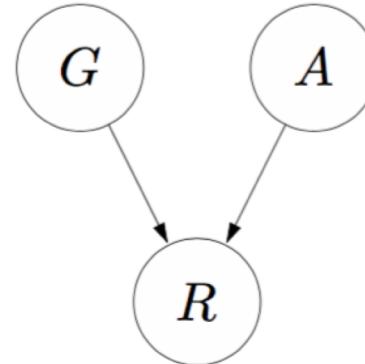
$g$	$r$	$p_R(r   g)$
d	4	2/3
d	5	1/3
c	1	1/2
c	5	1/2

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: Three Variables Example I

## Variables:

- Genre  $G \in \{\text{drama, comedy}\}$
- Won award  $A \in \{0, 1\}$
- Rating  $R \in \{1, 2, 3, 4, 5\}$



$$\mathbb{P}(G = g, A = a, R = r) = p_G(g)p_A(a)p_R(r | g, a)$$

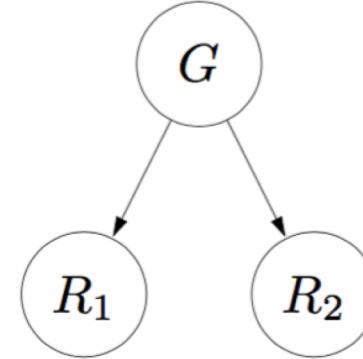
$$\mathcal{D}_{\text{train}} = \{(\text{d}, 0, 3), (\text{d}, 1, 5), (\text{c}, 0, 1), (\text{c}, 0, 5), (\text{c}, 1, 4)\}$$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: Three Variables Example II

## Variables:

- Genre  $G \in \{\text{drama, comedy}\}$
- Jim's rating  $R_1 \in \{1, 2, 3, 4, 5\}$
- Martha's rating  $R_2 \in \{1, 2, 3, 4, 5\}$



$$\mathbb{P}(G = g, R_1 = r_1, R_2 = r_2) = p_G(g)p_{R_1}(r_1 | g)p_{R_2}(r_2 | g)$$

$$\mathcal{D}_{\text{train}} = \{(d, 4, 5), (d, 4, 4), (d, 5, 3), (c, 1, 2), (c, 5, 4)\}$$

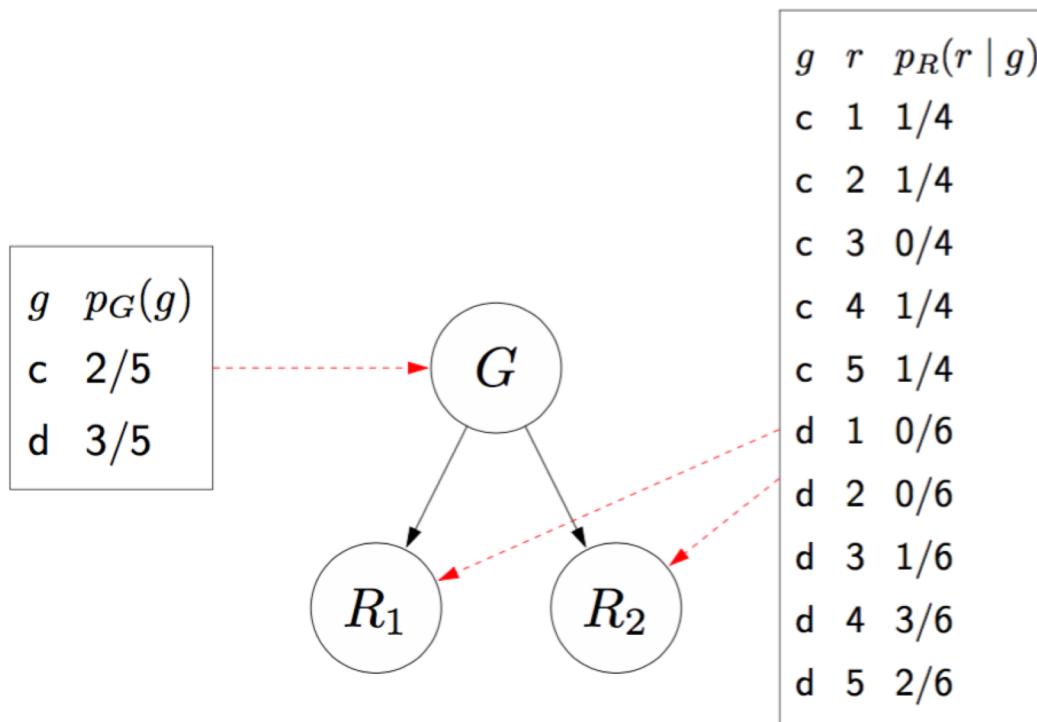
<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: Parameter Sharing



**Key idea: parameter sharing**

The local conditional distributions of different variables use the same parameters.



<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: Maximum Likelihood via Counting and Normalizing

**Input:** training examples  $\mathcal{D}_{\text{train}}$  of full assignments

**Output:** parameters  $\theta = \{p_d : d \in D\}$



**Algorithm: maximum likelihood for Bayesian networks**

**Count:**

For each  $x \in \mathcal{D}_{\text{train}}$ :

    For each variable  $x_i$ :

        Increment  $\text{count}_{d_i}(x_{\text{Parents}(i)}, x_i)$

**Normalize:**

For each  $d$  and local assignment  $x_{\text{Parents}(i)}$ :

    Set  $p_d(x_i | x_{\text{Parents}(i)}) \propto \text{count}_d(x_{\text{Parents}(i)}, x_i)$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: Maximum Likelihood via Counting and Normalizing

---

Maximum likelihood objective:

$$\max_{\theta} \prod_{x \in \mathcal{D}_{\text{train}}} \mathbb{P}(X = x; \theta)$$

Algorithm on previous slide exactly computes maximum likelihood parameters (closed form solution).

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Learning in DPGM: Maximum Likelihood via Counting and Normalizing

$$\mathcal{D}_{\text{train}} = \{(d, 4), (d, 5), (c, 5)\}$$

$$\max_{p_G(\cdot)} (p_G(d)p_G(d)p_G(c)) \max_{p_R(\cdot|c)} p_R(5 | c) \max_{p_R(\cdot|d)} (p_R(4 | d)p_R(5 | d))$$

- Key: decomposes into subproblems, one for each distribution  $d$  and assignment  $x_{\text{Parents}}$
- For each subproblem, solve in closed form (Lagrange multipliers for sum-to-1 constraint)

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Different Estimation/Learning Problems

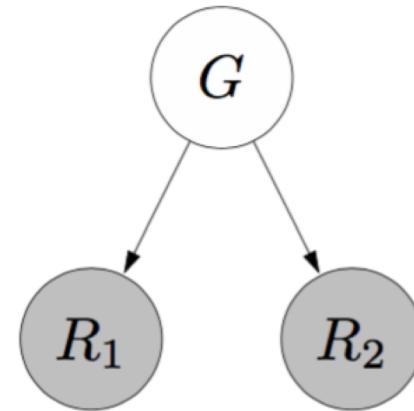
---

- What if we have missing data?

<b>Model</b>	DPGM	UPGM
<b>Data</b>	Complete	Incomplete
<b>Structure</b>	Known	Unknown
<b>Objective</b>	Generative	Discriminative

# Learning in DPGM: Latent Variables

---



What if we **don't observe** some of the variables?

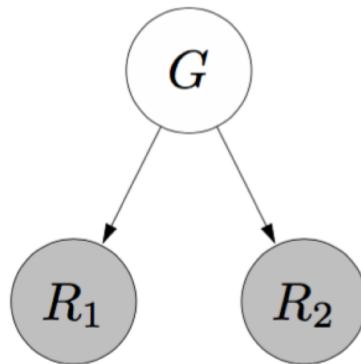
$$\mathcal{D}_{\text{train}} = \{(\textcolor{red}{?}, 4, 5), (\textcolor{red}{?}, 4, 4), (\textcolor{red}{?}, 5, 3), (\textcolor{red}{?}, 1, 2), (\textcolor{red}{?}, 5, 4)\}$$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# DPGM: Maximizing Marginal Likelihood

Variables:  $H$  is hidden,  $E = e$  is observed

Example:



$$H = G \quad E = (R_1, R_2) \quad e = (4, 5)$$
$$\theta = (p_G, p_R)$$

Maximum marginal likelihood objective:

$$\begin{aligned} & \max_{\theta} \prod_{e \in \mathcal{D}_{\text{train}}} \mathbb{P}(E = e; \theta) \\ &= \max_{\theta} \prod_{e \in \mathcal{D}_{\text{train}}} \sum_h \mathbb{P}(H = h, E = e; \theta) \end{aligned}$$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Expectation Maximization

Inspiration: K-means

Variables:  $H$  is hidden,  $E$  is observed (to be  $e$ )



## Algorithm: Expectation Maximization (EM)

E-step:

- Compute  $q(h) = \mathbb{P}(H = h | E = e; \theta)$  for each  $h$  (use any probabilistic inference algorithm)
- Create weighted points:  $(h, e)$  with weight  $q(h)$

M-step:

- Compute maximum likelihood (just count and normalize) to get  $\theta$

Repeat until convergence.

# EM: Revisiting K-Means

---

- EM tries to maximize marginal likelihood
- K-means
  - Is a special case of EM (for GMMs with variance tending to 0)
  - **Objective:** Estimate cluster centers

# EM: Revisiting K-Means

---

- EM tries to maximize marginal likelihood
- K-means
  - Is a special case of EM (for GMMs with variance tending to 0)
  - Objective: Estimate cluster centers
  - But don't know which points belong to which clusters
  - Take an alternating optimization approach
    - Find the best cluster assignment given current cluster centers
    - Find the best cluster centers given assignments

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# The Two Steps of EM

---

- E-step
  - Here, we don't know what the hidden variables are, so compute their distribution given the current parameters ( $P(H|E = e; \theta)$ )
  - Need inference! (BP/Gibbs MCMC)
  - This distribution provides a weight  $q(h)$  (temp variable in the EM algo) to every value  $H$  can take
- Conceptually, the E-step generates a set of weighted full realizations/configurations  $(h, e)$  with weights  $q(h)$

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

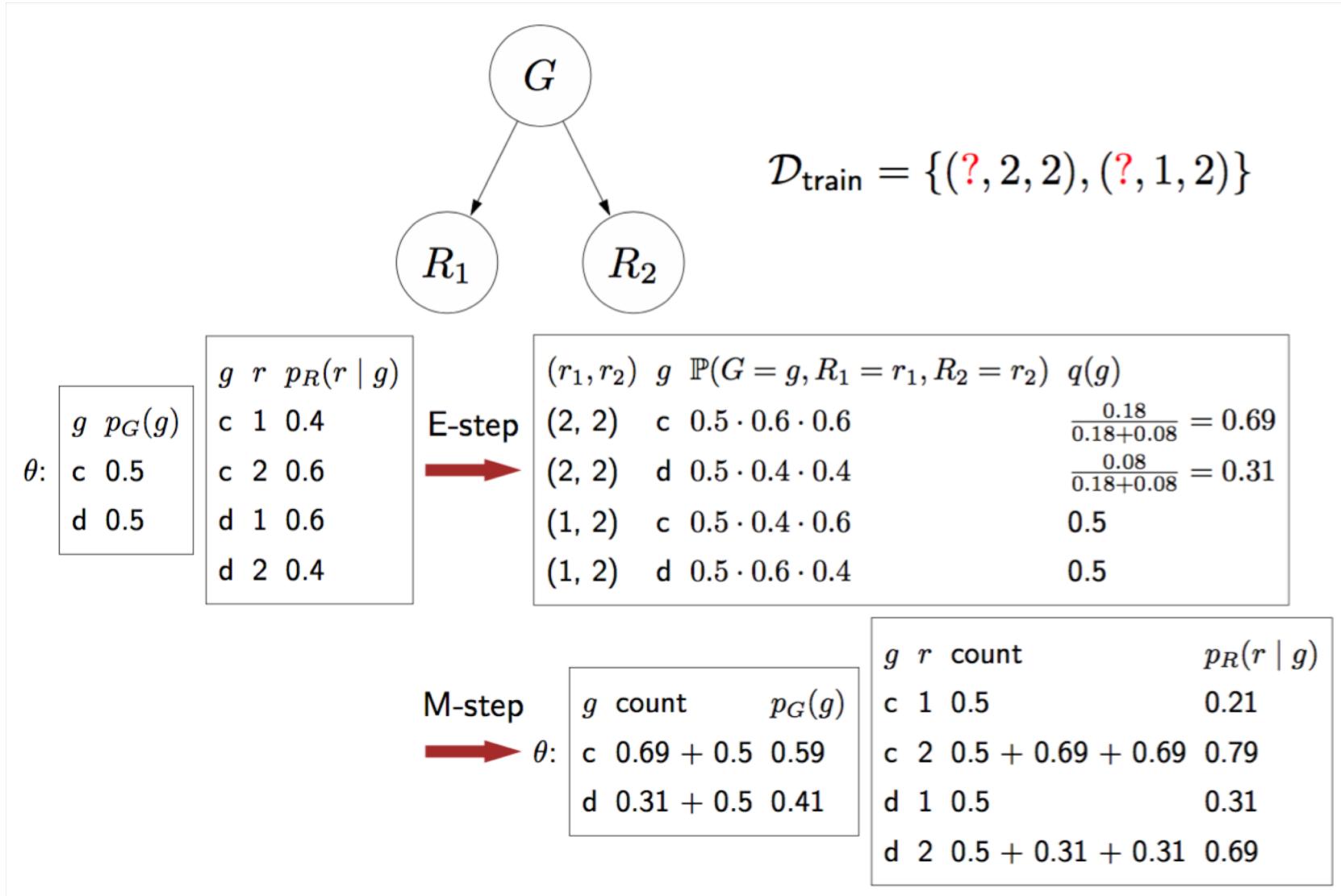
# The Two Steps of EM

---

- M-step
  - Just do MLE (i.e., counting and normalizing) to re-estimate parameters
- If we repeat E-step and M-step again and again, eventually we will converge to a local optima of parameters

<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# EM: Example



<sup>1</sup>Reference: Percy Liang, CS221 (2015)

# Different Estimation/Learning Problems

---

- What if the structure is unknown?

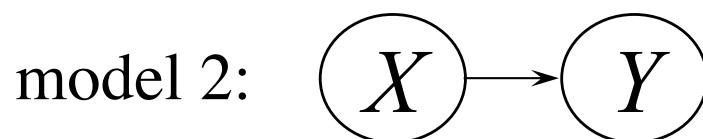
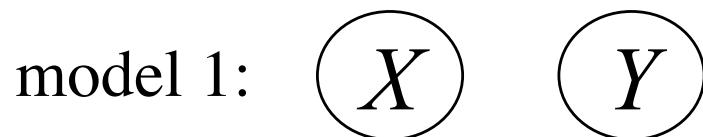
<b>Model</b>	DPGM	UPGM
<b>Data</b>	Complete	Incomplete
<b>Structure</b>	Known	Unknown
<b>Objective</b>	Generative	Discriminative

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Learning Structure: Bayesian Approach

---

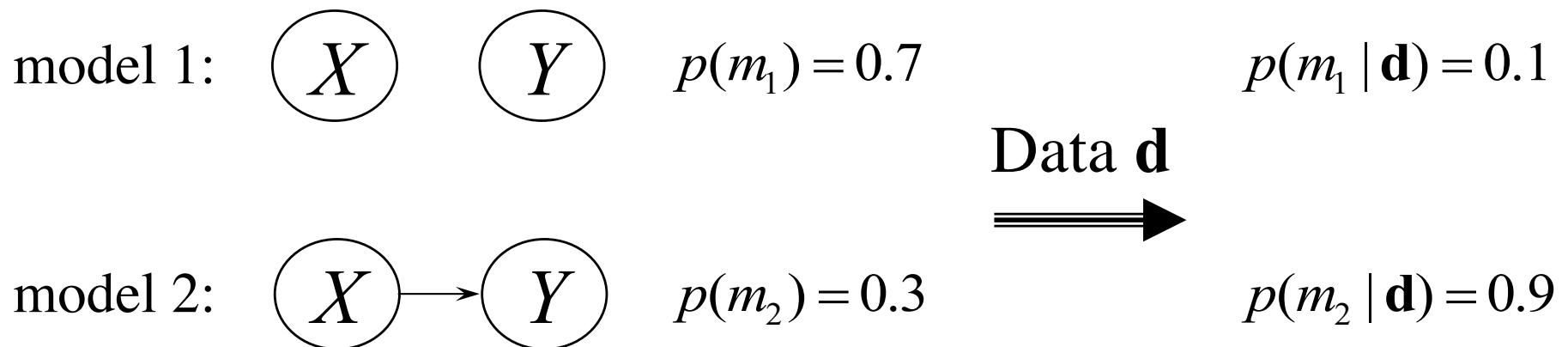
- Given data, which model is correct?



<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Learning Structure: Bayesian Approach

- Given data, which model is ~~correct?~~ more likely?



- Can do model averaging
- Can do model selection to pick a model that is
  - tractable, understandable, explainable

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Learning Structure: Model Scoring

- Use Baye's theorem to score a model

Given data  $\mathbf{d}$ :

$$\text{model score} \longrightarrow p(m | \mathbf{d}) \propto p(m) \underbrace{p(\mathbf{d} | m)}$$

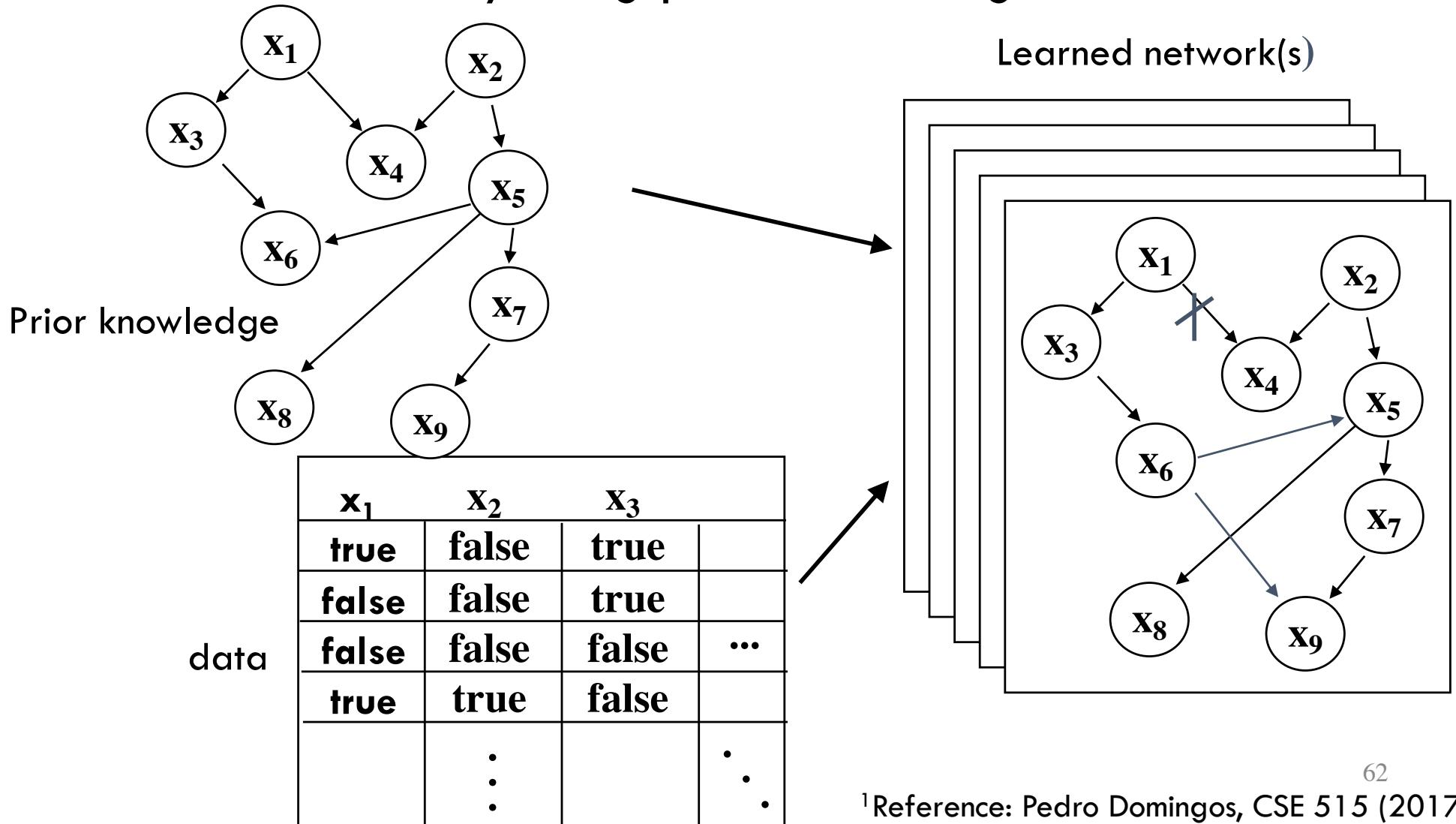
"marginal likelihood"

$$p(\mathbf{d} | m) = \int p(\mathbf{d} | \theta, m) p(\theta | m) d\theta$$

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Combined Learning

- Although structure learning is hard in general, still useful to do it by using prior knowledge and data



<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Different Estimation/Learning Problems

---

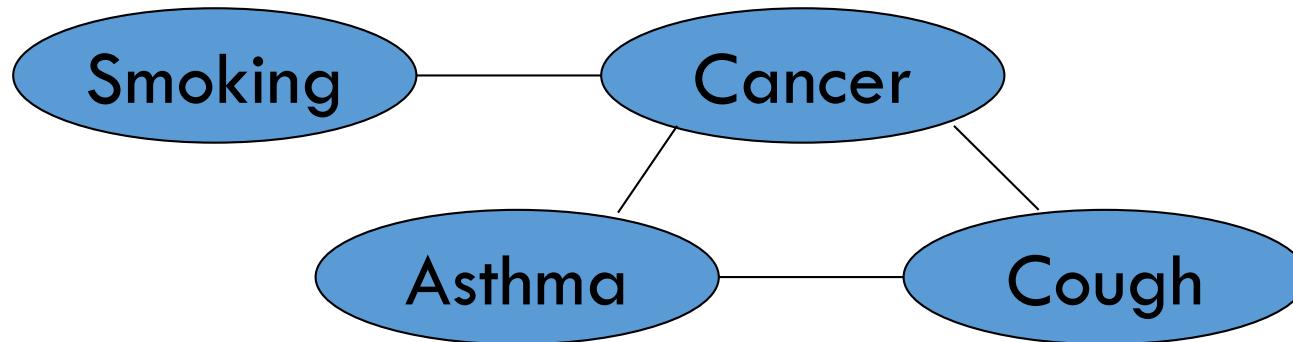
- There are many variants

<b>Model</b>	DPGM	UPGM
<b>Data</b>	Complete	Incomplete
<b>Structure</b>	Known	Unknown
<b>Objective</b>	Generative	Discriminative

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Learning in UPGM

---



Potential functions defined over cliques

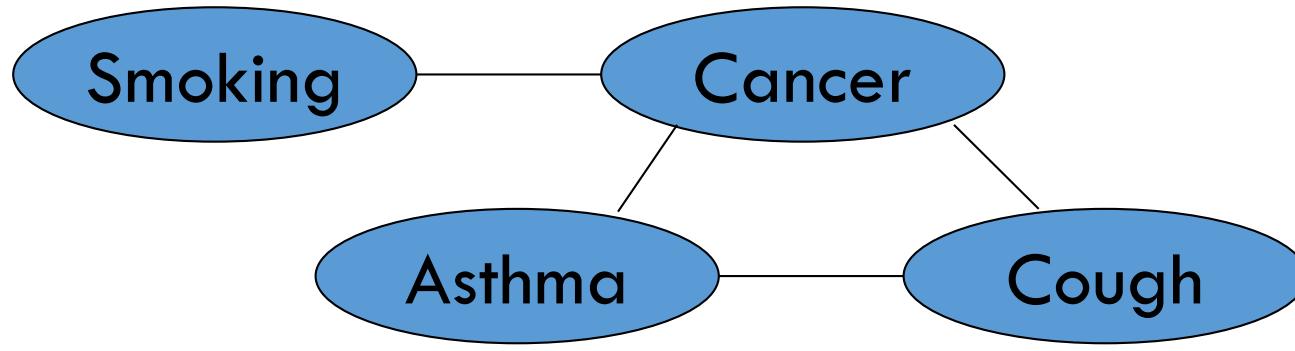
$$P(x) = \frac{1}{Z} \prod_c \Phi_c(x_c)$$

$$Z = \sum_x \prod_c \Phi_c(x_c)$$

Smoking	Cancer	$\Phi(S,C)$
False	False	4.5
False	True	4.5
True	False	2.7
True	True	4.5

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Learning in UPGM



Can be thought in terms of a log-linear representation

$$P(x) = \frac{1}{Z} \exp\left( \sum_i w_i f_i(x) \right)$$

Weight of Feature  $i$     Feature  $i$

$$f_1(\text{Smoking}, \text{Cancer}) = \begin{cases} 1 & \text{if } \neg \text{Smoking} \vee \text{Cancer} \\ 0 & \text{otherwise} \end{cases}$$

$$w_1 = 0.51$$

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Learning in UPGM: Generative

---

- Maximize likelihood or posterior probability
- Numerical optimization (gradient or 2<sup>nd</sup> order)

$$\frac{\partial}{\partial w_i} \log P_w(x) = n_i(x) - E_w[n_i(x)]$$

No. of times feature  $i$  is true in data

Expected no. times feature  $i$  is true according to model

- Requires inference at each step (slow!)

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Learning in UPGM: Pseudo-likelihood

---

$$PL(x) \equiv \prod_i P(x_i | \text{neighbors}(x_i))$$

- Likelihood of each variable given its neighbors in the data
- Does not require inference at each step
- Consistent estimator
- Widely used in vision, spatial statistics, etc.
- But PL parameters may not work well for long inference chains

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Different Estimation/Learning Problems

---

- There are many variants

<b>Model</b>	DPGM	UPGM
<b>Data</b>	Complete	Incomplete
<b>Structure</b>	Known	Unknown
<b>Objective</b>	Generative	Discriminative

# Learning in UPGM: Discriminative

---

- This is related to Conditional Random Fields (CRFs)
- Maximize conditional likelihood of query ( $y$ ) given evidence ( $x$ )

$$\frac{\partial}{\partial w_i} \log P_w(y | x) = n_i(x, y) - E_w[n_i(x, y)]$$

No. of true values of feature  $i$  in data

Expected no. of true values according to model

- Inference is easier, but still hard

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Different Estimation/Learning Problems

---

- There are many variants

<b>Model</b>	DPGM	UPGM
<b>Data</b>	Complete	Incomplete
<b>Structure</b>	Known	Unknown
<b>Objective</b>	Generative	Discriminative

# Learning in UPGM: Missing Data

---

- Gradient of likelihood is now difference of expectations

$$\frac{\partial}{\partial w_i} \log P_w(x) = E_w[n_i(y | x)] - E_w[n_i(x, y)]$$

Exp. no. true values given observed data

$x$ : Observed

$y$ : Missing

Expected no. true values given no data

- Can use gradient descent or EM

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Learning Summary

---

- We looked at the following problems
  - Learning DPGMs with complete data and known structure
    - MLE via counting and normalizing
  - Learning DPGMs with incomplete data and known structure
    - EM
  - Learning DPGM structure
  - Learning UPGMs in a generative setting
  - Learning UPGM in a discriminative setting

# Learning Summary

---

- There are many other variants
- Some of these tasks necessarily rely on heuristics
- Many ways have been proposed in research, and as practitioners, we have to pick and choose.

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

---

---

# Questions?

# Summary

---

- We discussed some of the applications where they have been successfully applied
- We looked at parameter and structure estimation of these graphical models
- Bottom line: When there is structure in the inputs and outputs of a ML pipeline, consider DPGMs/UPGMs
  - An unified way of thinking about supervised and unsupervised learning

---

---

# Appendix

# Sample Exam Questions

---

- In which settings would one use MLE and EM for learning in graphical models? Give examples.
- How is the graph structure learned? Can it be specified as prior information?
- Mention 3 applications of graphical models and specify their descriptions. Explain how learning happens in one of these models.

Which is computationally more expensive for Bayesian networks?

probabilistic inference given the parameters

learning the parameters given fully labeled data

# Gibbs Sampling when Observations/Evidence are Given

“State” of network = current assignment to all variables.

Generate next state by sampling one variable given Markov blanket  
Sample each variable in turn, keeping evidence fixed

```
function GIBBS-SAMPLING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
    local variables:  $\mathbf{N}[X]$ , a vector of counts over  $X$ , initially zero
                     $\mathbf{Z}$ , the nonevidence variables in  $bn$ 
                     $\mathbf{x}$ , the current state of the network, initially copied from  $e$ 
    initialize  $\mathbf{x}$  with random values for the variables in  $\mathbf{Y}$ 
    for  $j = 1$  to  $N$  do
        for each  $Z_i$  in  $\mathbf{Z}$  do
            sample the value of  $Z_i$  in  $\mathbf{x}$  from  $P(Z_i|mb(Z_i))$ 
            given the values of  $MB(Z_i)$  in  $\mathbf{x}$ 
             $\mathbf{N}[x] \leftarrow \mathbf{N}[x] + 1$  where  $x$  is the value of  $X$  in  $\mathbf{x}$ 
    return NORMALIZE( $\mathbf{N}[X]$ )
```

Can also choose a variable to sample at random each time

<sup>1</sup>Reference: Pedro Domingos, CSE 515 (2017)

# Additional Applications: Naïve Bayes Spam Filter

---

- **Key assumption**

Words occur independently of each other given the label of the document

$$p(w_1, \dots, w_n | \text{spam}) = \prod_{i=1}^n p(w_i | \text{spam})$$

- **Spam classification via Bayes Rule**

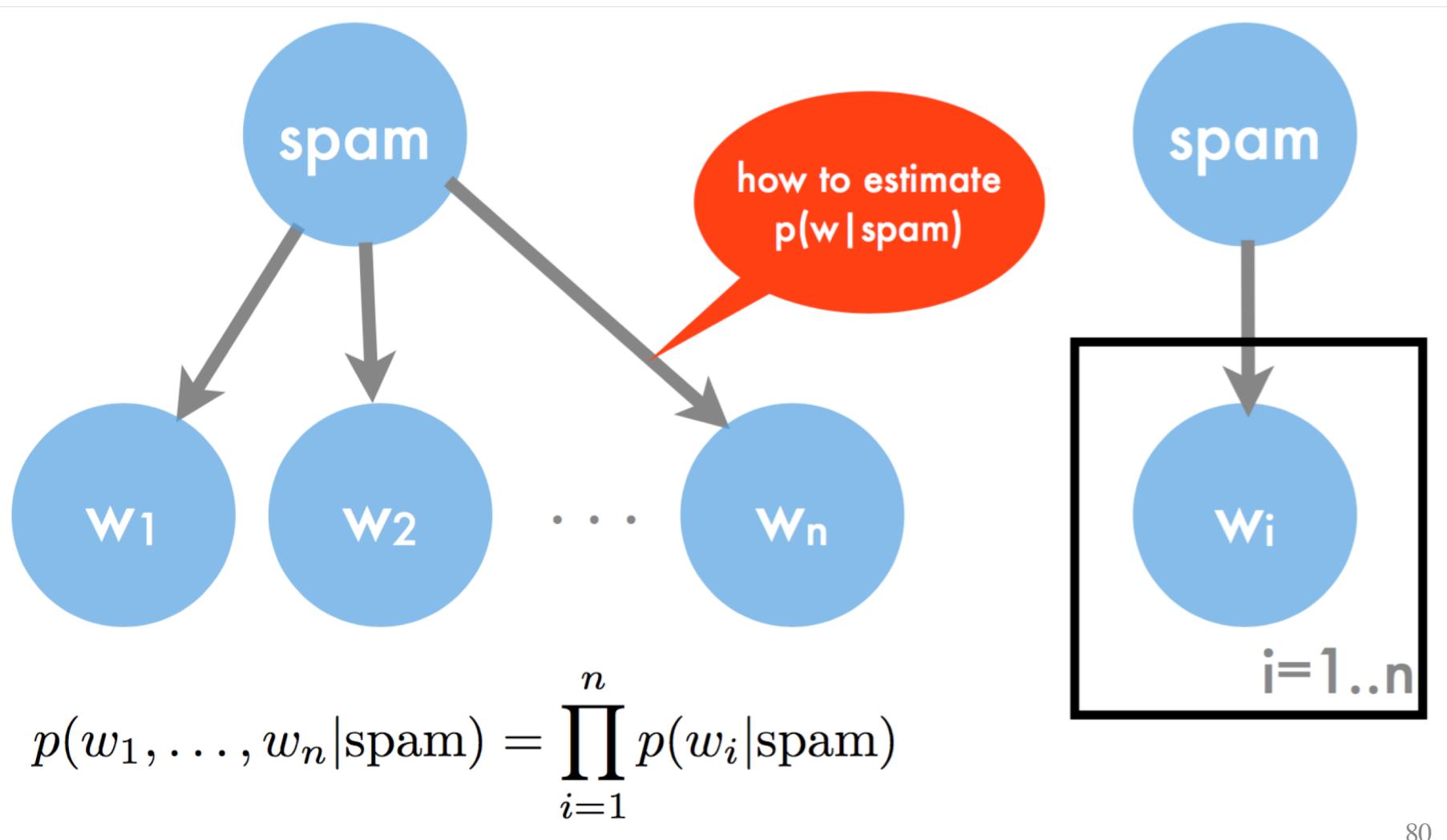
$$p(\text{spam} | w_1, \dots, w_n) \propto p(\text{spam}) \prod_{i=1}^n p(w_i | \text{spam})$$

- **Parameter estimation**

Compute spam probability and word distributions for spam and ham

<sup>1</sup>Reference: Alex Smola (2011)

# Additional Applications: Naïve Bayes Spam Filter



<sup>1</sup>Reference: Alex Smola (2011)

# Additional Applications: Naïve Bayes Spam Filter

- Two classes (spam/ham)
- Binary features (e.g. presence of \$\$\$, viagra)
- Simplistic Algorithm
  - Count occurrences of feature for spam/ham
  - Count number of spam/ham mails

feature probability

spam probability

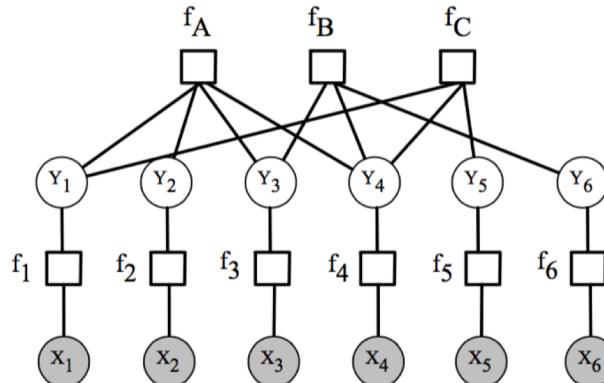
$$p(x_i = \text{TRUE}|y) = \frac{n(i, y)}{n(y)} \text{ and } p(y) = \frac{n(y)}{n}$$

$$p(y|x) \propto \frac{n(y)}{n} \prod_{i:x_i=\text{TRUE}} \frac{n(i, y)}{n(y)} \prod_{i:x_i=\text{FALSE}} \frac{n(y) - n(i, y)}{n(y)}$$

<sup>1</sup>Reference: Alex Smola (2011)

# Additional Applications: MAP Problem in Low Density Parity Check Codes

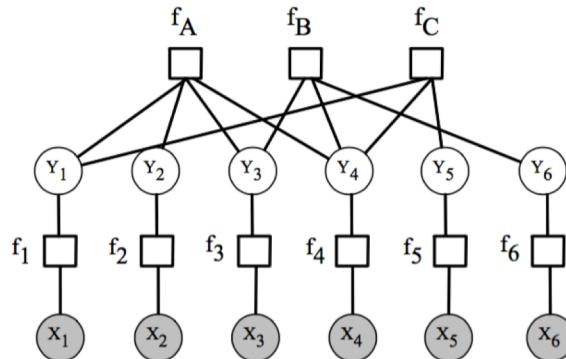
- Error correcting codes for transmitting a message over a noisy channel (invented by Gallager in the 1960's, then re-discovered in 1996)



- Each of the top row factors enforce that its variables have even parity:  
 $f_A(Y_1, Y_2, Y_3, Y_4) = 1$  if  $Y_1 \otimes Y_2 \otimes Y_3 \otimes Y_4 = 0$ , and 0 otherwise
- Thus, the only assignments  $\mathbf{Y}$  with non-zero probability are the following (called **codewords**): *3 bits encoded using 6 bits*  
000000, 011001, 110010, 101011, 111100, 100101, 001110, 010111
- $f_i(Y_i, X_i) = p(X_i | Y_i)$ , the likelihood of a bit flip according to noise model

<sup>1</sup>Reference: David Sontag (2013)

# Additional Applications: MAP Problem in Low Density Parity Check Codes



- The *decoding* problem for LDPCs is to find

$$\operatorname{argmax}_{\mathbf{y}} p(\mathbf{y} \mid \mathbf{x})$$

This is called the **maximum a posteriori** (MAP) assignment

- Since  $Z$  and  $p(\mathbf{x})$  are constants with respect to the choice of  $\mathbf{y}$ , can equivalently solve (taking the log of  $p(\mathbf{y}, \mathbf{x})$ ):

$$\operatorname{argmax}_{\mathbf{y}} \sum_{c \in C} \theta_c(\mathbf{x}_c),$$

where  $\theta_c(\mathbf{x}_c) = \log \phi_c(\mathbf{x}_c)$

<sup>1</sup>Reference: David Sontag (2013)