

Case Study 4: Collaborative Filtering

Matrix Factorization and Probabilistic LFs for Network Modeling

Machine Learning for Big Data
CSE547/STAT548, University of Washington

Emily Fox

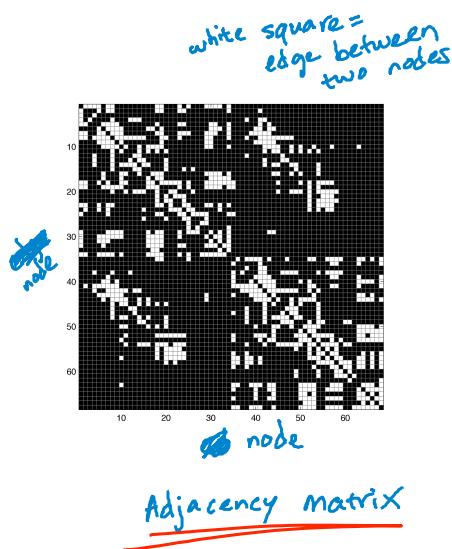
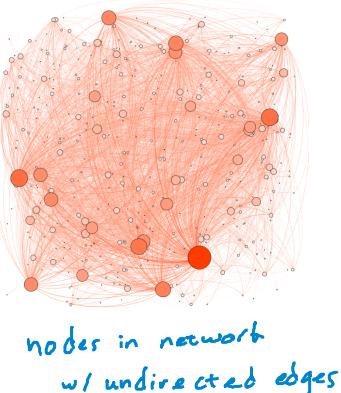
May 21st, 2015

©Emily Fox 2015

1

Network Data

■ Structure of network data



©Emily Fox 2015

2

Properties of Data Source

- Similarities to Netflix data:
 - Matrix-valued data (adj. matrix)
 - High-dimensional many nodes
 - Sparse few links between nodes (e.g. ppl in a social network)
- Differences
 - ✗ □ Square ← same indices for rows + columns
 - ✗ □ Binary ← yes/no for link (other ext. possible... multiple)

If undirected, then matrix is Symmetric



©Emily Fox 2015

3

Matrix Factorization for Network Data

- Vanilla matrix factorization approach:

In undirected case, just introduce node (e.g. user) factors L_u
 $r_{uv} \approx L_u \cdot L_v$ ← edge between users $u + v$

In directed case, just introduce sender factors L_u and receiver factors \tilde{L}_v
 $r_{uv} \approx L_u \cdot \tilde{L}_v$ ← edge from user u to user v
- What to return for link prediction?

Is r_{uv} binary? $L_u \in \mathbb{R}^k \rightarrow$ no
 Many options, but can return top k
 $r_{uv_1}, r_{uv_2}, \dots, r_{uv_k}$ (just threshold rule)
- Slightly fancier:
 More appropriate to have $r_{uv} \in [0, 1]$
 use $r_{uv} = \sigma(L_u \cdot L_v)$ $\sigma = \text{logistic fcn}$

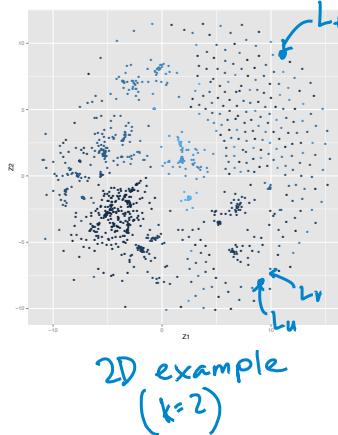
©Emily Fox 2015

4

Probabilistic Latent Space Models

- Assume features (covariates) of the user X_u or relationship X_{uv} \leftarrow could just be $X_u X_v$
- Each user has a “position” in a k -dimensional latent space
unobserved (learn) $L_u \in \mathbb{R}^k$, just as in matrix fact.
- Probability of link:

$$\begin{aligned} \text{log odds } p(r_{uv}=1 | L_u, L_v, x_{uv}, \beta) &= \log \frac{p(r_{uv}=1 | L_u, L_v, x_{uv}, \beta)}{p(r_{uv}=0 | L_u, L_v, x_{uv}, \beta)} \\ &= \beta_0 + \beta^T x_{uv} - \|L_u - L_v\| \\ &\quad \text{OR} \\ &= \beta_0 + \beta^T x_{uv} + \|L_u^T L_v\| \end{aligned}$$



©Emily Fox 2015

5

Probabilistic Latent Space Models

- Probability of link:

$$\text{log odds } p(r_{uv} = 1 | L_u, L_v, x_{uv}, \beta) = \beta_0 + \beta^T x_{uv} - \|L_u - L_v\|$$

prob. link is high for L_u close to L_v

$$\text{log odds } p(r_{uv} = 1 | L_u, L_v, x_{uv}, \beta) = \beta_0 + \beta^T x_{uv} + \|L_u^T L_v\|$$

can modify as $\frac{\|L_u^T L_v\|}{\|L_v\|}$

prob. of link is high if β bt $L_u + L_v$ small

- Bayesian approach:

- Place prior on user factors and regression coefficients
- Place hyperprior on user factor hyperparameters

- Many other options and extensions (e.g., can use GMM for $L_u \rightarrow$ clustering of users in the latent space)

©Emily Fox 2015

6

What you need to know...

- Representation of network data as a matrix
 - Adjacency matrix
- Similarities and differences between adjacency matrices and general matrix-valued data
- Matrix factorization approaches for network data
 - Just use standard MF and threshold output
 - Introduce link functions to constrain predicted values
- Probabilistic latent space models
 - Model link probabilities using distance between latent factors

©Emily Fox 2015

7

Case Study 5: Mixed Membership Modeling

Clustering Documents Revisited, Latent Dirichlet Allocation

Machine Learning for Big Data
CSE547/STAT548, University of Washington
Emily Fox
May 21st, 2015

©Emily Fox 2015

8

Document Retrieval

■ **Goal:** Retrieve documents of interest

■ **Challenges:**

- Tons of articles out there
- How should we measure similarity?



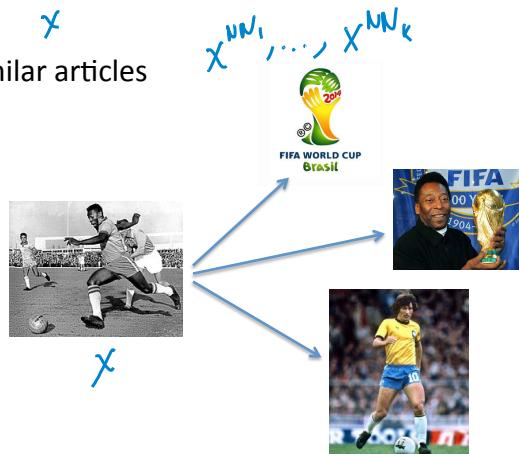
©Emily Fox 2015

9

Task 1: Find Similar Documents

■ **First considered:**

- **Input:** Query article x
- **Output:** Set of k similar articles



©Emily Fox 2015

10

Task 2: Cluster Documents

■ Then examined:

- Cluster documents based on topic



©Emily Fox 2015

11

Document Representation

■ Bag of words model



$w_i = 3$
 $\Rightarrow i^{\text{th}}$ word
in this doc
is "hat"

previously ↗ vector fan
of word counts
 $X = \begin{bmatrix} & \\ & \end{bmatrix}$
(e.g. tf-idf)

performed operations on
this vector

representation of a doc

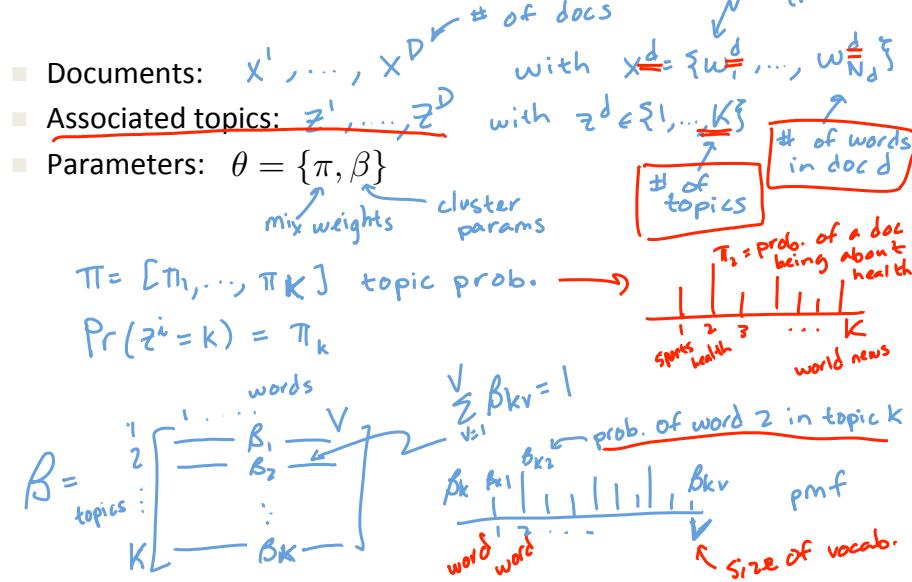
representation of a doc

$X = \{w_1, \dots, w_N\}$ # of words in doc
unordered set of N words in doc.
 $w_i \in V$ (vocab)

©Emily Fox 2015

12

A Generative Model



©Emily Fox 2015

13

A Generative Model

- Documents: x^1, \dots, x^D
- Associated topics: z^1, \dots, z^D
- Parameters: $\theta = \{\pi, \beta\}$
- Generative model:

Sample topic:

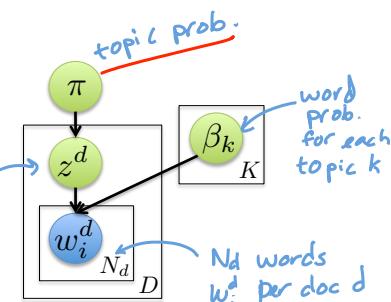
$$z^d \sim \pi$$

sample words:

$$w_i^d | z^d \sim \beta_{z^d} \quad i=1, \dots, N_d$$

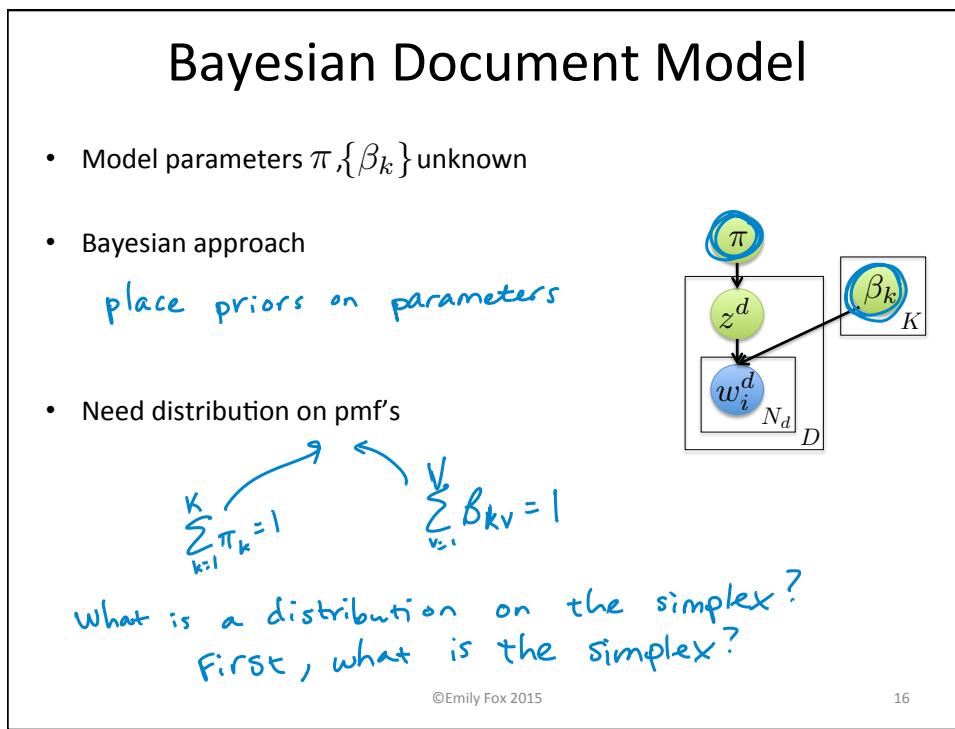
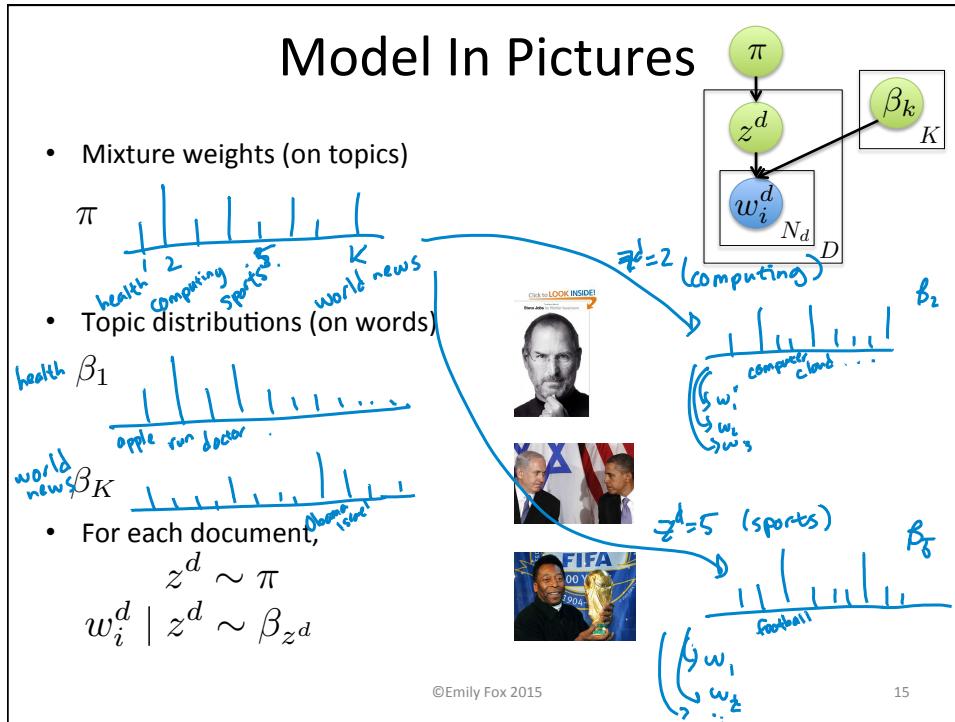
topic of doc d

Given topic $z^d = k$ for doc d , draw each word from β_k



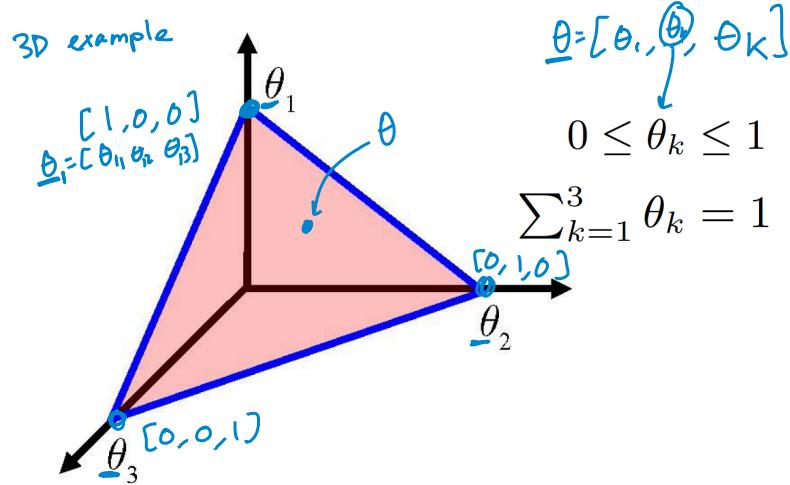
©Emily Fox 2015

14



The Simplex in 3D

- The simplex defines the hyperplane of vectors that sum to 1

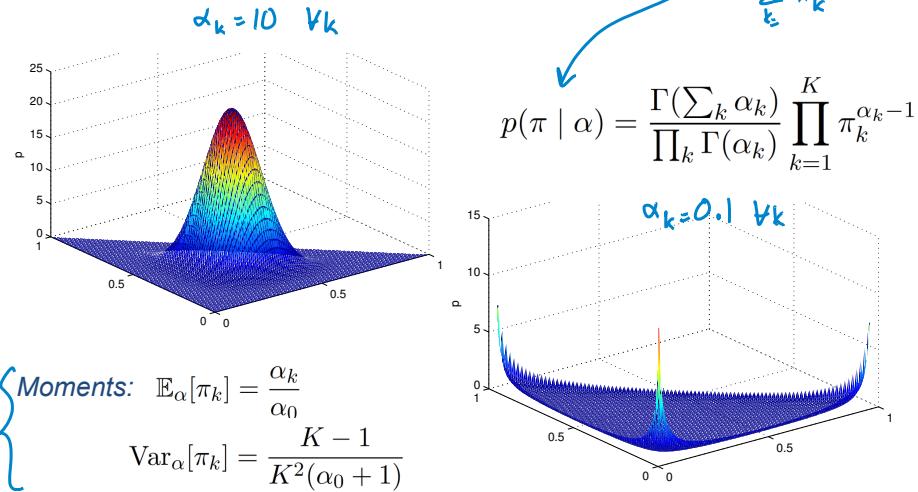


©Emily Fox 2015

17

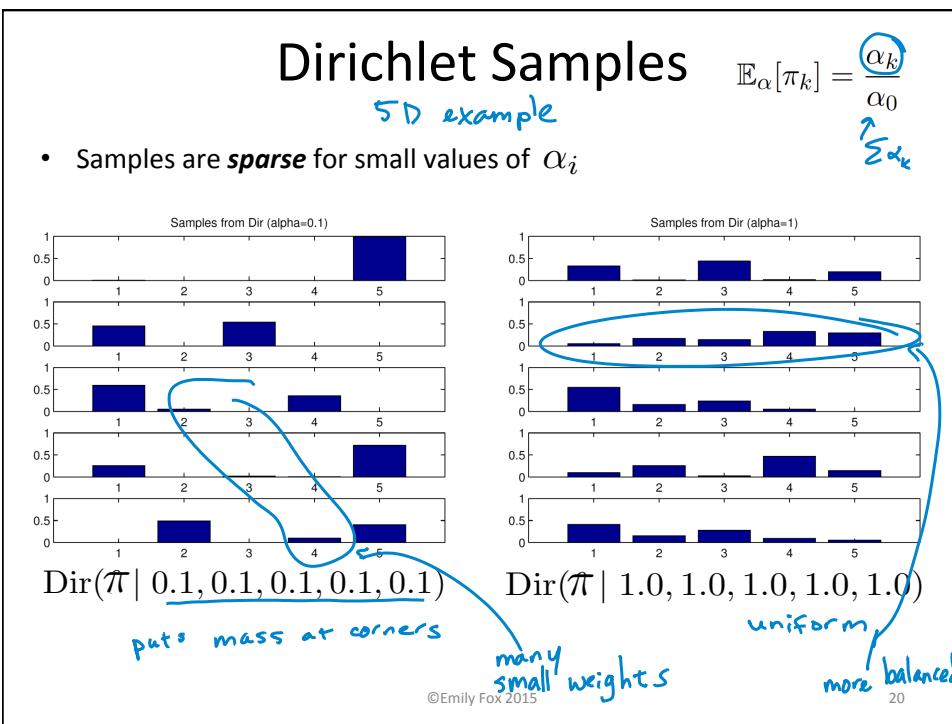
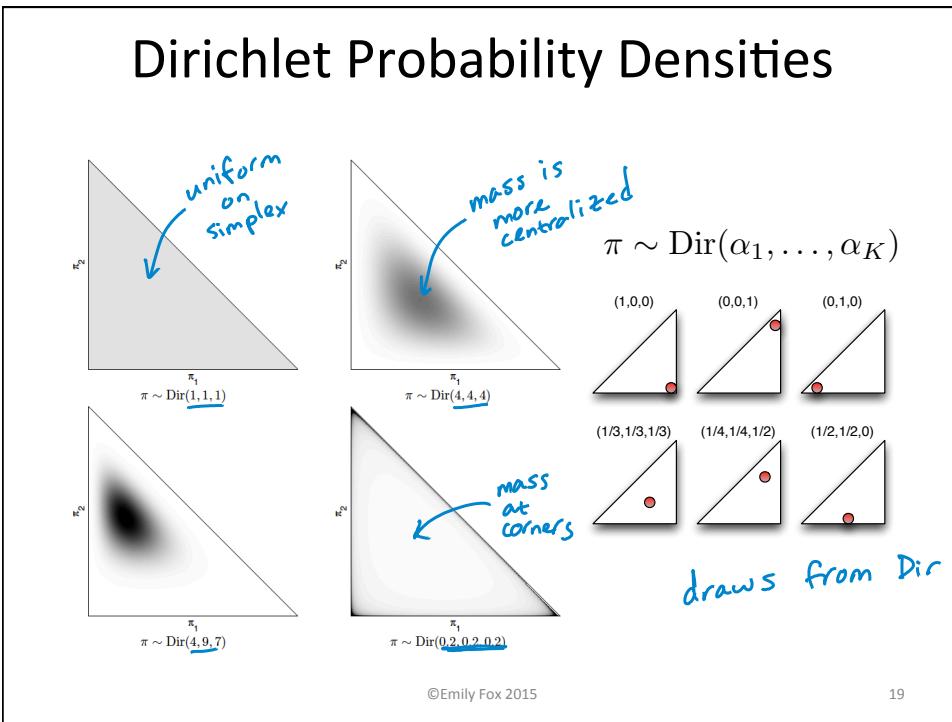
Dirichlet Distributions

- The Dirichlet distribution is defined on the simplex



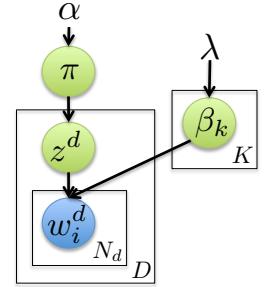
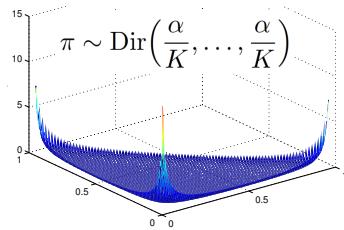
©Emily Fox 2015

18



Model Summary

- Prior on model parameters
 - E.g., symmetric Dirichlet for π



- Dirichlet prior for topic parameters $\beta_k \sim \text{Dir}\left(\frac{\lambda}{V}, \dots, \frac{\lambda}{V}\right)$
- Sample observations as

$$\begin{aligned} z^d &\sim \pi & d = 1, \dots, D \\ w_i^d \mid z^d &\sim \underline{\beta_{z^d}} & i = 1, \dots, N_d \end{aligned}$$

©Emily Fox 2015

21

Posterior Inference via Sampling

- Iterate between sampling

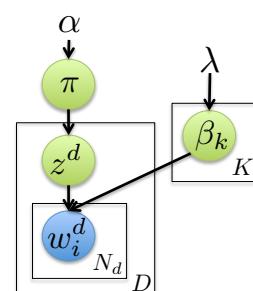
$$\pi \sim p(\pi \mid \{z^d\}, \{\beta_k\}, \{w_i^d\})$$

For $k=1, \dots, K$

$$\beta_k \sim p(\beta_k \mid \pi, \{z^d\}, \{w_i^d\})$$

For $d=1, \dots, D$

$$z^d \sim p(z^d \mid \pi, \{\beta_k\}, \{w_i^d\})$$



- What form do these complete conditionals take?

- First a look at statements of conditional independence in directed graphical models

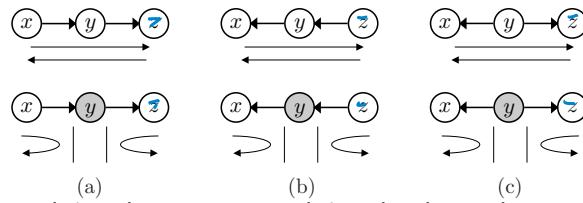
©Emily Fox 2015

22

Conditional Independence in Bayes Nets

$$p(\text{joint}) = \prod_i p(x_i | \text{parents}(x_i))$$

- Consider 4 different junction configurations



- Conditional versus unconditional independence:

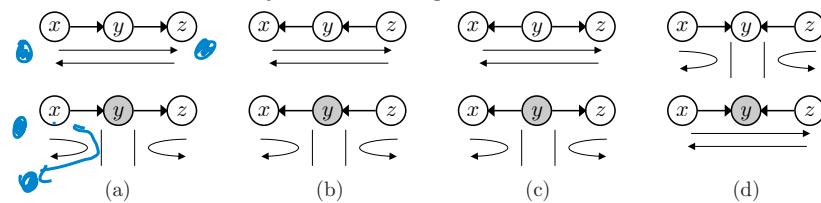
$$\sum_y p(x, y, z) = p(x)p(z)p(y|x, z) \Rightarrow p(x, z) = p(x)p(z)$$

$$p(x, z|y) \propto p(x, y, z) = p(x)p(z)p(y|x, z) \quad x \perp\!\!\!\perp z \mid y$$

"explaining away": $x = \text{earthquake}$ $z = \text{burglar}$ $y = \text{alarm}$
 but if alarm=1, then ©Emily Fox 2015 $\begin{matrix} \nearrow \text{ind. a priori} \\ \searrow \text{increase in prob}(z) \Rightarrow \text{decrease in prob}(x) \end{matrix}$

Bayes Ball Algorithm

- Consider 4 different junction configurations



- Bayes ball algorithm

Start ball at one end or the other.

If ball passes to a node (straight) then,
not cond./marg. independent.

If ball bounces back (wall + curved arrows),
 then nodes are cond/marg ind.

Markov Blanket

- A node is conditionally independent of all other nodes in the graph given its Markov blanket

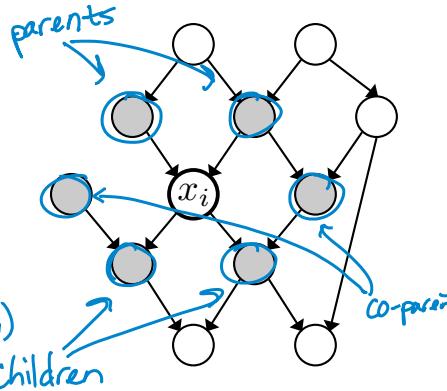
*Markov
blanket* = - all parents
of
 x_i
- all children
- all coparents

- Gibbs sampling iterates between full conditionals

$$x_i \sim p(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N)$$

→ simplify to

$$x_i \sim p(x_i | MB(x_i))$$



©Emily Fox 2015

25

What you need to know...

- Bayesian specification of document clustering model
- Rules of conditional and unconditional independence in directed graphical models (Bayes nets)
 - Bayes' ball
 - Markov blanket

©Emily Fox 2015

26

Reading for *Next Lecture*

- **Mixed Membership Models: KM Sec. 27.3**
 - Basic LDA:
[Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." the Journal of machine Learning research 3 \(2003\): 993-1022.](#)
 - Introduction:
[Blei, David M. "Probabilistic topic models." Communications of the ACM, vol. 55, no. 4 \(2012\): 77-84.](#)
 - Sampling:
[Griffith, Thomas L. and Mark Steyvers. "Finding scientific topics." Proceedings of the National Academy of Sciences of the United States of America, Volume: 101, Supplement: 1 \(2004\): Pages: 5228-5235](#)