

Case Study 4: Collaborative Filtering

GraphLab Review,

ALS, SGD, Gibbs Sampling for MF in GraphLab

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

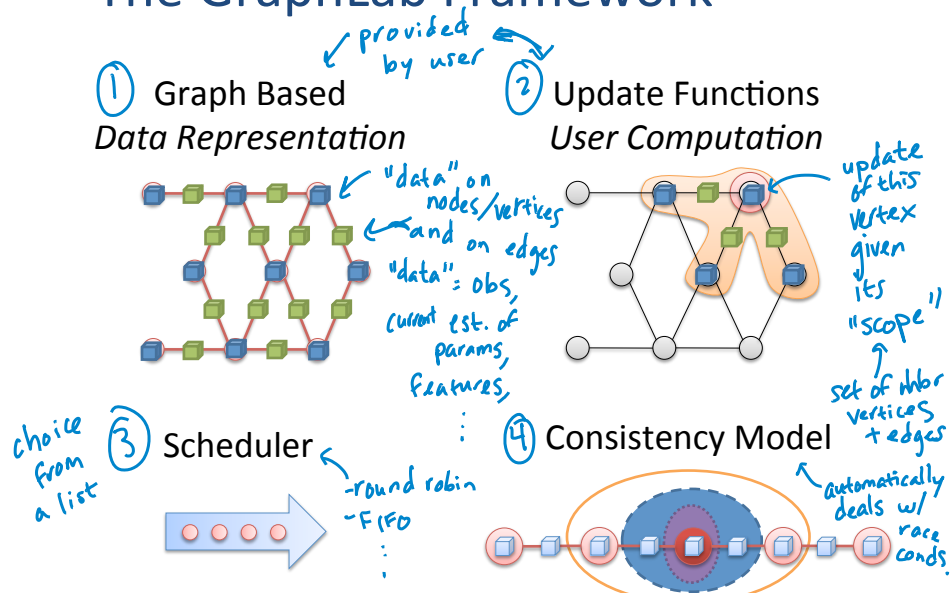
Emily Fox

February 27th, 2014

©Emily Fox 2014

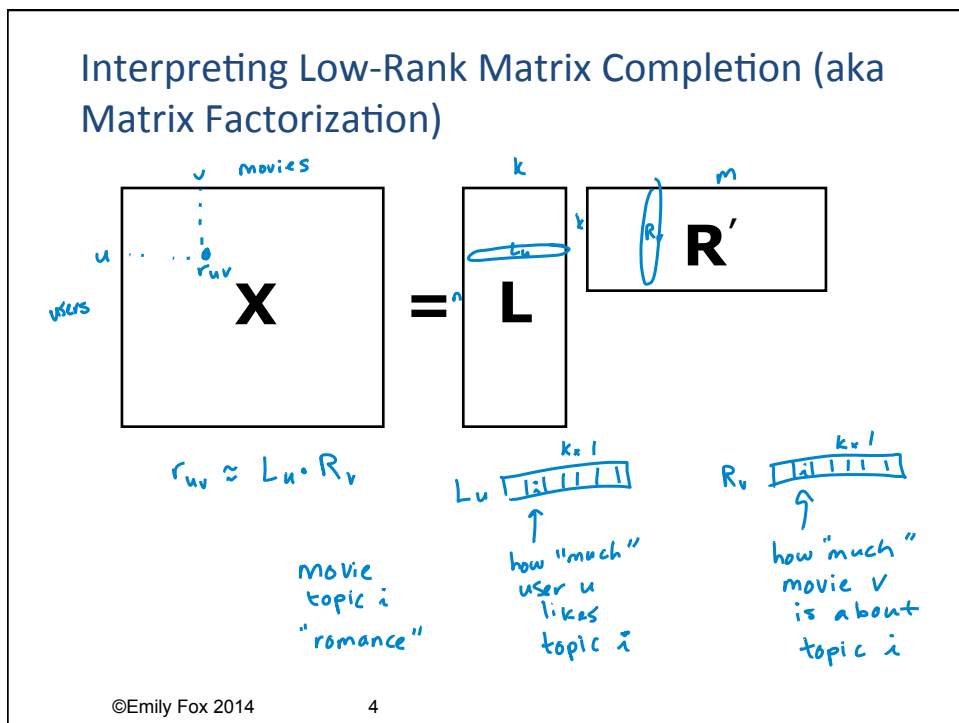
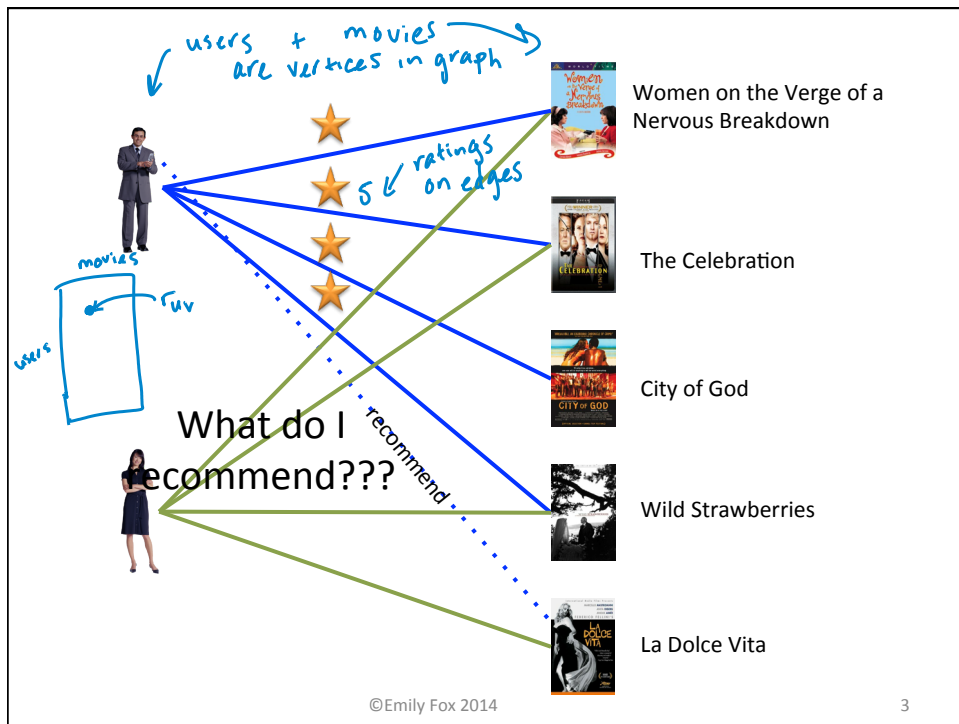
1

The GraphLab Framework

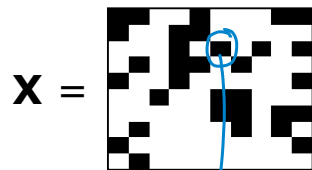


©Emily Fox 2014

2

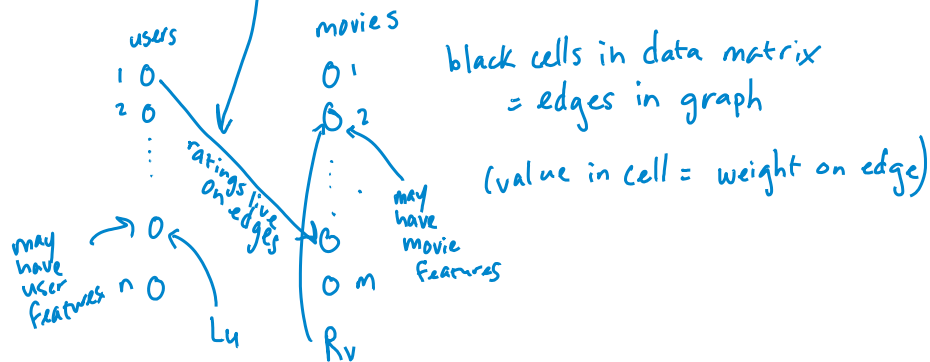


Matrix Completion as a Graph



$X =$

X_{ij} known for black cells
 X_{ij} unknown for white cells
 Rows index users
 Columns index movies



©Emily Fox 2014

5

Coordinate Descent for Matrix Factorization: Alternating Least-Squares

$$\min_{L, R} \sum_{(u,v): r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\| + \lambda_v \|R\|$$

model w/o user/movie features

- Fix movie factors, optimize for user factors

□ Independent least-squares over users

$$\min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\|$$

nhbrs of user u (movies rated by user u)

- Fix user factors, optimize for movie factors

□ Independent least-squares over movies

$$\min_{R_v} \sum_{u \in U_v} (L_u \cdot R_v - r_{uv})^2 + \lambda_v \|R\|$$

nhbrs of movie v

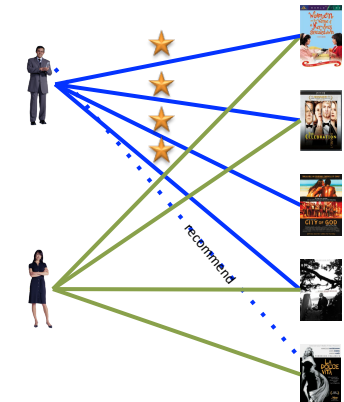
- System may be underdetermined: use regularization

- Converges to local optima

©Emily Fox 2014

6

Alternating Least Squares Update Function



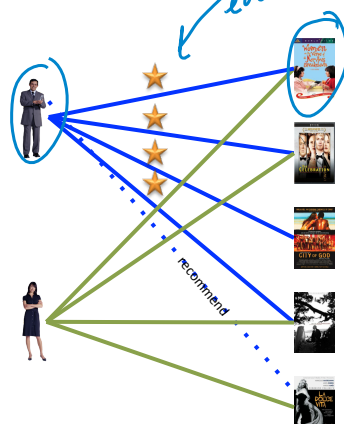
$$\min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 + \lambda \|L_u\|^2$$

$$\min_{R_v} \sum_{u \in U_v} (L_u \cdot R_v - r_{uv})^2 + \lambda \|R_v\|^2$$

update(i, scope) {
 read current factors for nhbrs
 build e.g. $X = \begin{bmatrix} \text{---} \\ -R_v- \\ \text{---} \\ \vdots \end{bmatrix}$ } all movies rated by user i
 read all ratings on edges
 build $y = \begin{bmatrix} r_{i,v} \\ \vdots \end{bmatrix} \leftarrow v \in V_u$
 Solve local regression problem
 e.g. for L_2 reg,
 $L_u = (X^T X + \lambda_u I)^{-1} X^T y$
 }

©Emily Fox 2014 7

SGD for Matrix Factorization in GraphLab



$$\epsilon_t = L_u^{(t)} \cdot R_v^{(t)} - r_{uv}$$

$$\begin{bmatrix} L_u^{(t+1)} \\ R_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_u) L_u^{(t)} - \eta_t \epsilon_t R_v^{(t)} \\ (1 - \eta_t \lambda_v) R_v^{(t)} - \eta_t \epsilon_t L_u^{(t)} \end{bmatrix}$$

GraphLab operates on vertices
 Update(i, scope) {
 Perform SGD update
 for each neighbor of i
 (for every edge connected to i)
 }

©Emily Fox 2014 8

Bayesian PMF Example

* Full Bayesian approach
place priors on ϕ as well!

- Latent user and movie factors:

$$L_u \sim N(\mu_u, \Sigma_u) \quad u=1, \dots, n$$

$$R_v \sim N(\mu_v, \Sigma_v) \quad v=1, \dots, m$$

- Observations $r_{uv} \sim N(L_u^T R_v, \sigma_r^2)$

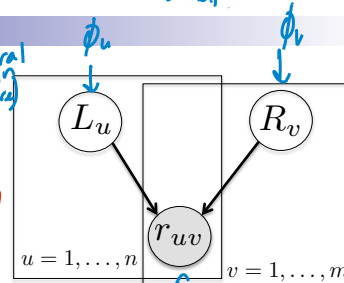
- Hyperparameters:

$$\phi = \{\underbrace{\mu_u, \Sigma_u}_{\phi_u}, \underbrace{\mu_v, \Sigma_v}_{\phi_v}, \underbrace{\sigma_r^2}_{\phi_r}\}$$

- Want to predict new movie rating:

$$p(r_{uv}^* | X, \phi) = \int p(r_{uv}^* | L_u, R_v) p(L, R | X, \phi) dL dR$$

\uparrow new rating \uparrow obs. ratings \uparrow new user/movie combo \uparrow posterior given obs. so far



©Emily Fox 2014

9

Bayesian PMF Gibbs Sampler

- Outline of Bayesian PMF sampler

1. Init $L^{(1)}, R^{(1)}$

2. For $k=1, \dots, \text{Niter}$

(i) Sample hyperparams $\phi^{(k)} = \{\phi_u^{(k)}, \phi_v^{(k)}, \phi_r^{(k)}\}$

(ii) For each user $u=1, \dots, n$ sample in parallel

$$L_u^{(k+1)} \sim P(L_u | X, R^{(k)}, \phi^{(k)})$$

← form on next slide

(iii) For each movie $v=1, \dots, m$ sample in parallel

$$R_v^{(k+1)} \sim P(R_v | X, L^{(k+1)}, \phi^{(k)})$$

Very similar to ideas of ALS (systematically)

©Emily Fox 2014

10

Bayesian PMF Example

- For user u :

$$p(L_u | X, R, \phi_u) \propto p(L_u | \phi_u) \prod_{v \in V_u} p(r_{uv} | L_u, R_v, \phi_r)$$

$$\propto N(L_u | \mu_u, \Sigma_u) \prod_{v \in V_u} N(r_{uv} | L_u \cdot R_v, \sigma_r^2)$$

$$= N(L_u | \tilde{\mu}_u, \tilde{\Sigma}_u) \quad \leftarrow \text{via conjugacy}$$

where $\tilde{\Sigma}_u^{-1} = \Sigma_u^{-1} + \sigma_r^{-2} \sum_{v \in V_u} R_v R_v^T$ (sum over all nhbrs of user u)

$\tilde{\mu}_u = \tilde{\Sigma}_u (\sigma_r^{-2} \sum_{v \in V_u} r_{uv} R_v + \Sigma_u^{-1} \mu_u)$ (vertex weight \leftarrow edge weights) (posterior is in the same family as prior)

- Symmetrically for R_v conditioned on L (breaks down over movies)
- Luckily, we can use this to get our desired posterior samples

©Emily Fox 2014

11

PMF Gibbs Sampling in GraphLab

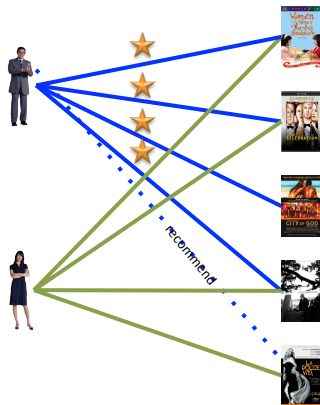
$$p(L_u | X, R, \phi_u) = N(\tilde{\mu}_u, \tilde{\Sigma}_u) \quad \tilde{\Sigma}_u^{-1} = \Sigma_u^{-1} + \sigma_r^{-2} \sum_{v \in V_u} R_v R_v^T \quad \tilde{\mu}_u = \tilde{\Sigma}_u \left(\sigma_r^{-2} \sum_{v \in V_u} r_{uv} R_v + \Sigma_u^{-1} \mu_u \right)$$

update(i , scope) {
e.g. user i (update to L_i)

read current factors for nhbrs R_j
read ratings on edges r_{ij}
set $\tilde{\Sigma}_i^{-1} = \Sigma_i^{-1} + \sigma_r^{-2} \sum_{j \in \text{nhbrs}(i)} R_j R_j^T$
fixed at vertex i

set $\tilde{\mu}_i = \tilde{\Sigma}_i / \sigma_r^{-2} \sum_{j \in \text{nhbrs}(i)} r_{ij} R_j + \Sigma_i^{-1} \mu_i$

sample $L_i \sim N(\tilde{\mu}_i, \tilde{\Sigma}_i)$



©Emily Fox 2014

12



Release 2.2 available now

<http://graphlab.org>

Documentation... Code... Tutorials... (more on the way)

GraphChi 0.1 available now

<http://graphchi.org>

What you need to know...

- Data-parallel versus graph-parallel computation
- Bulk synchronous processing versus asynchronous processing
- GraphLab system for graph-parallel computation
 - Data representation
 - Update functions
 - Scheduling
 - Consistency model
- ALS, SGD and Gibbs for matrix factorization/PMF in GraphLab

*collaborative
filtering*

Reading

- Papers under “Case Study IV: **Parallel Learning with GraphLab**”
- Optional:
 - Parallel Splash BP
<http://www.ml.cmu.edu/research/dap-papers/dap-gonzalez.pdf>

Acknowledgements

- Slides based on Carlos Guestrin’s GraphLab talk

Case Study 5: Mixed Membership Modeling

Clustering Documents Revisited, Latent Dirichlet Allocation

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

Emily Fox

February 27th, 2014

©Emily Fox 2014

17

Document Retrieval

- **Goal:** Retrieve documents of interest

- **Challenges:**

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



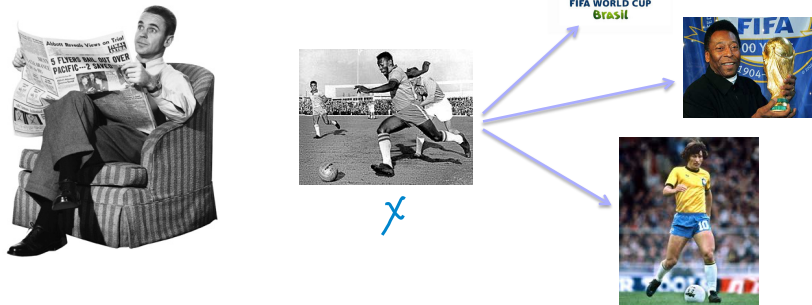
©Emily Fox 2014

18

Task 1: Find Similar Documents

■ First considered:

- Input: Query article x
- Output: Set of k similar articles $x^{MN_1}, \dots, x^{MN_k}$



©Emily Fox 2014

19

Task 2: Cluster Documents

■ Then examined:

- Cluster documents based on topic



©Emily Fox 2014

20

Document Representation

■ Bag of words model



document d

$d=3$
 $w_2=3$
 \Rightarrow 3rd word
 in doc d
 is 'hat'

cat	← word 1	2
dog		2
hat		3
car		...
...		...

previously: $x^d = \begin{bmatrix} \vdots \end{bmatrix}$ in NN search or some clustering alg.
 vector fcn of word counts (e.g. tf-idf)

representation of doc d
 performed operations on this vector

now: $x^d = \{w_1^d, \dots, w_{N_d}^d\}$ "bag of words"
 unordered set of N_d word w with $w_i^d \in V$ vocab.
 indices

©Emily Fox 2014

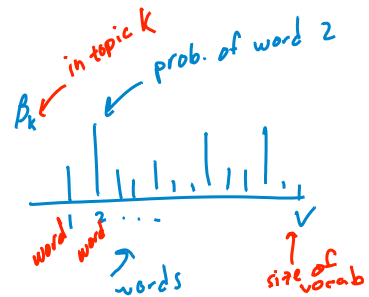
21

A Generative Model

- Documents: x^1, \dots, x^D with $x^d = \{w_1^d, \dots, w_{N_d}^d\}$
- Associated topics: z^1, \dots, z^D with $z^d \in \{1, \dots, K\}$
- Parameters: $\theta = \{\pi, \beta\}$ topic indicator for doc 1 ↑ # topics

as before $\left\{ \begin{array}{l} \pi = [\pi_1, \dots, \pi_K] \text{ topic probabilities} \\ \Pr(z^d = k) = \pi_k \end{array} \right. \leftarrow \Pr(z^d = k | \pi) = \pi_k$

$$\beta = \begin{matrix} & 1 & 2 & \dots & V \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ K \end{matrix} & \left[\begin{array}{cccc} \beta_1 & & & \\ & \beta_2 & & \\ & & \ddots & \\ & & & \beta_K \end{array} \right] \end{matrix}$$

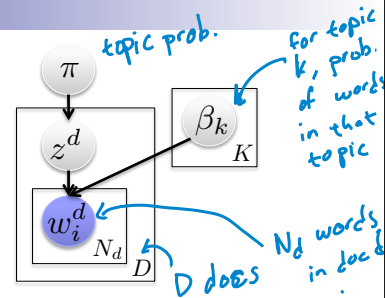


©Emily Fox 2014

22

A Generative Model

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



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

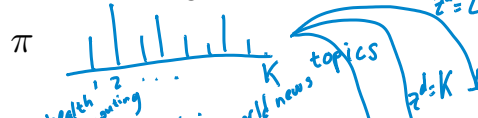
Given topic $z^d = k$ for doc d , draw each word from $\beta_k \leftarrow$ word prob. for topic k

©Emily Fox 2014

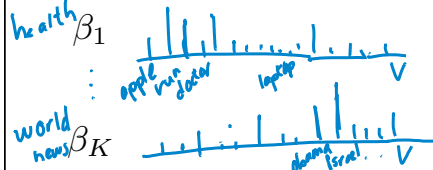
23

Model In Pictures

- Mixture weights (on topics)

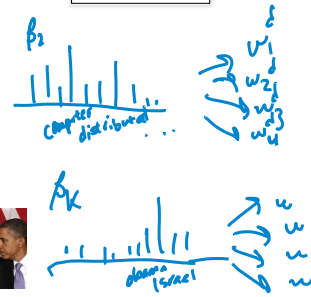
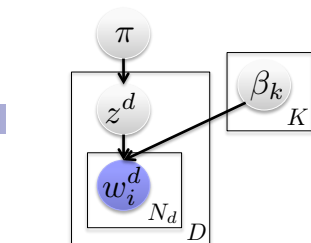


- Topic distributions (on words)



- For each document,

$z^d \sim \pi$
 $w_i^d | z^d \sim \beta_{z^d}$



©Emily Fox 2014

24

Bayesian Document Model

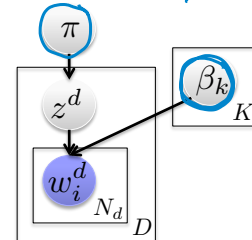
- Model parameters $\pi, \{\beta_k\}$ unknown

← can use EM as in case study 2

- Bayesian approach

place priors on parameters

- Need distribution on pmf's



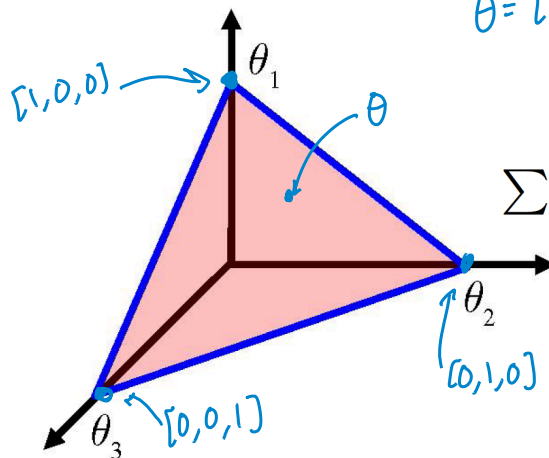
$\sum_{k=1}^K \pi_k = 1$ $\sum_{k=1}^K \beta_{kv} = 1$ ← π, β_k live on the simplex
 ① What is the simplex?
 ② What is a distribution on the simplex?

©Emily Fox 2014

25

The Simplex in 3D

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



$\theta = [\theta_1, \dots, \theta_k]$
 e.g. $k=3$ here

$$0 \leq \theta_k \leq 1$$

$$\sum_{k=1}^3 \theta_k = 1$$

©Emily Fox 2014

26

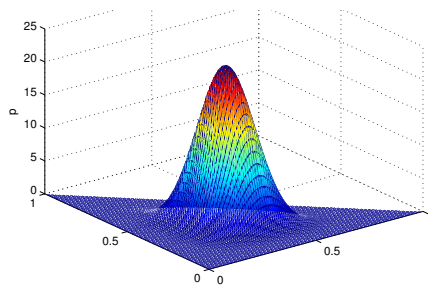
Dirichlet Distributions

- The Dirichlet distribution is defined on the simplex

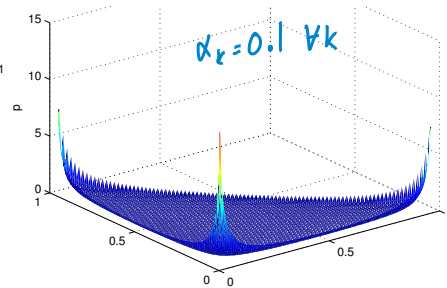
$$\alpha_k = 10 \quad \forall k$$

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$$

$$\Rightarrow \sum \pi_k = 1 \text{ and } \pi_k \geq 0 \quad \forall k$$



$$p(\pi | \alpha) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k - 1}$$



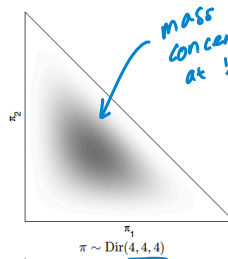
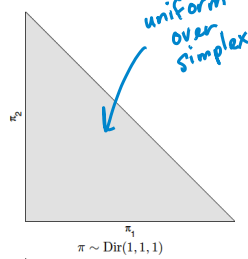
Moments: $\mathbb{E}_\alpha[\pi_k] = \frac{\alpha_k}{\alpha_0}$

$$\text{Var}_\alpha[\pi_k] = \frac{K-1}{K^2(\alpha_0+1)}$$

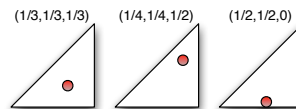
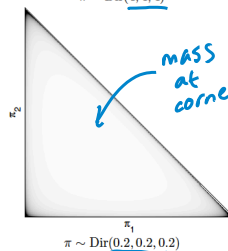
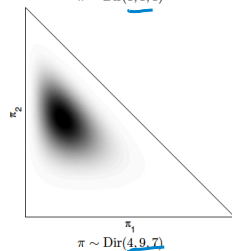
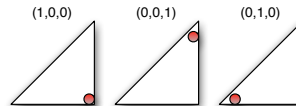
©Emily Fox 2014

27

Dirichlet Probability Densities



$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$$



draws from Dir

©Emily Fox 2014

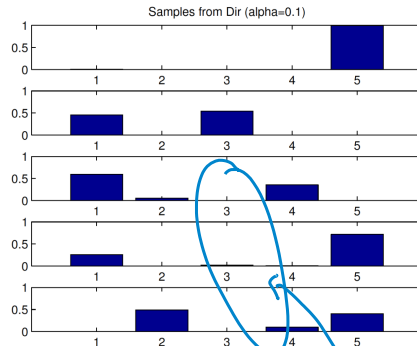
28

Dirichlet Samples

$$\mathbb{E}_{\alpha}[\pi_k] = \frac{\alpha_k}{\alpha_0}$$

- Samples are **sparse** for small values of α_i

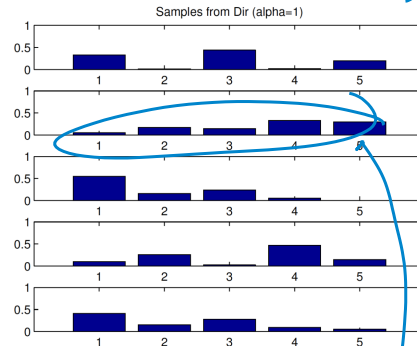
(5D example)



Dir(π | 0.1, 0.1, 0.1, 0.1, 0.1)

puts mass at corners

many small weights



Dir(π | 1.0, 1.0, 1.0, 1.0, 1.0)

much more uniform

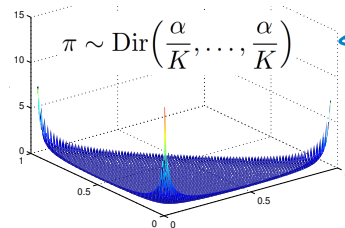
©Emily Fox 2014

29

Model Summary

- Prior on model parameters

- E.g., symmetric Dirichlet for π

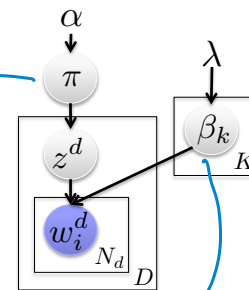


- Dirichlet prior for topic parameters $\beta_k \sim \text{Dir}(\lambda_1, \dots, \lambda_v)$ $k=1, \dots, K$

- Sample observations as

$$z^d \sim \pi \quad d=1, \dots, D$$

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

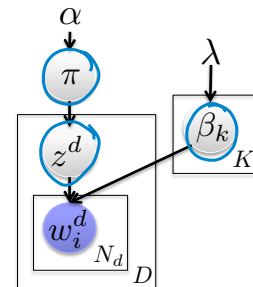


©Emily Fox 2014

30

Posterior Inference via Sampling

- Iterate between sampling *actual obs.*
 $\pi \sim p(\pi | \{z^d\}, \{\beta_k\}, \{w_i^d\})$
 For $k=1, \dots, K$
 $\beta_k \sim p(\beta_k | \pi, \{z^d\}, \{\beta_j, j \neq k\}, \{w_i^d\})$
 For $d=1, \dots, D$
 $z^d \sim p(z^d | \pi, \{z^i, i \neq d\}, \{\beta_k\}, \{w_i^d\})$



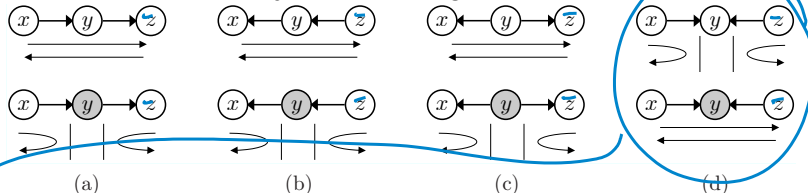
- What form do these complete conditionals take?
 - First a look at statements of conditional independence in directed graphical models

©Emily Fox 2014

31

Conditional Independence in Bayes Nets

- Consider 4 different function configurations



- Conditional versus unconditional independence:

$$p(x, y, z) = p(x)p(z)p(y|x, z) \stackrel{\text{int. over } y}{\Rightarrow} p(x, z) = p(x)p(z) \Rightarrow x \perp\!\!\!\perp z$$

$$p(x, z|y) \propto p(x, y, z) = p(x)p(z)p(y|x, z) \neq p(x|y)p(z|y) \leftarrow x \not\perp\!\!\!\perp z | y$$

"explaining away": $x = \text{earthquake}$, $z = \text{burglar}$, $y = \text{car alarm}$
ind. a priori
 If alarm ($y=1$), an increase in earthquake $p(x|y)$, means $p(z|y)$ lower

©Emily Fox 2014

32