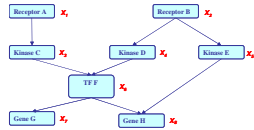


Approximate Inference: Markov Chain Monte Carlo

Probabilistic Graphical Models (10-708)

Lecture 19, Nov 26, 2007



Eric Xing

Reading: J-Chap. 1, KF-Chap. 11

1

Markov chain Monte Carlo (MCMC)



- Importance sampling does not scale well to high dimensions.
- Rao-Blackwellisation not always possible.
- MCMC is an alternative.
- Construct a Markov chain whose stationary distribution is the target density $= \mathcal{P}(X|e)$.
- Run for T samples (burn-in time) until the chain converges/mixes/reaches stationary distribution.
- Then collect M (correlated) samples x_m .
- Key issues:
 - Designing proposals so that the chain mixes rapidly.
 - Diagnosing convergence.

Markov Chains



- **Definition:**

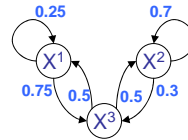
- Given an n-dimensional state space
- Random vector $\mathbf{X} = (X_1, \dots, X_n)$
- $\mathbf{x}^{(t)} = \mathbf{x}$ at time-step t
- $\mathbf{x}^{(t)}$ transitions to $\mathbf{x}^{(t+1)}$ with prob
 $P(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)}) = T(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}) = T(\mathbf{x}^{(t)} \rightarrow \mathbf{x}^{(t+1)})$

- **Homogenous:** chain determined by state $\mathbf{x}^{(0)}$, fixed **transition kernel** T (rows sum to 1)

- **Equilibrium:** $\pi(\mathbf{x})$ is a **stationary (equilibrium) distribution** if
 $\pi(\mathbf{x}') = \sum_{\mathbf{x}} \pi(\mathbf{x}) T(\mathbf{x} \rightarrow \mathbf{x}')$.

i.e., is a left eigenvector of the transition matrix $\pi^T T = \pi^T$.

$$(0.2 \ 0.5 \ 0.3) = (0.2 \ 0.5 \ 0.3) \begin{pmatrix} 0.25 & 0 & 0.75 \\ 0 & 0.7 & 0.3 \\ 0.5 & 0.5 & 0 \end{pmatrix}$$



Eric Xing

3

Markov Chains



- An MC is **irreducible** if transition graph connected
- An MC is **aperiodic** if it is not trapped in cycles
- An MC is **ergodic** (regular) if you can get from state x to x' in a finite number of steps.
- **Detailed balance:** $\text{prob}(\mathbf{x}^{(t)} \rightarrow \mathbf{x}^{(t-1)}) = \text{prob}(\mathbf{x}^{(t-1)} \rightarrow \mathbf{x}^{(t)})$

$$p(\mathbf{x}^{(t)}) T(\mathbf{x}^{(t-1)} | \mathbf{x}^{(t)}) = p(\mathbf{x}^{(t-1)}) T(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$$

summing over $\mathbf{x}^{(t-1)}$

$$p(\mathbf{x}^{(t)}) = \sum_{\mathbf{x}^{(t-1)}} p(\mathbf{x}^{(t-1)}) T(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)})$$

- Detailed bal \rightarrow stationary dist exists

Eric Xing

4

Metropolis-Hastings



- Treat the target distribution as stationary distribution
- Sample from an easier proposal distribution, followed by an acceptance test
- This induces a transition matrix that satisfies detailed balance

- MH proposes moves according to $Q(x'|x)$ and accepts samples with probability $A(x'|x)$.
- The induced transition matrix is $T(x \rightarrow x') = Q(x'|x)A(x'|x)$
- Detailed balance means

$$\pi(x)Q(x'|x)A(x'|x) = \pi(x')Q(x|x')A(x|x')$$

- Hence the acceptance ratio is

$$A(x'|x) = \min\left(1, \frac{\pi(x')Q(x|x')}{\pi(x)Q(x'|x)}\right)$$

Eric Xing

5

Metropolis-Hastings



1. Initialize $x^{(0)}$
 2. While not mixing // burn-in
 - $x = x^{(t)}$
 - $t += 1$,
 - sample $u \sim \text{Unif}(0,1)$
 - sample $x^* \sim Q(x^*|x)$
 - if $u < A(x^*|x) = \min\left(1, \frac{\pi(x^*)Q(x|x^*)}{\pi(x)Q(x^*|x)}\right)$
 - $x^{(t)} = x^*$ // transition
 - else
 - $x^{(t)} = x$ // stay in current state
- } Function
Draw sample $x(t)$
- Reset $t=0$, for $t=1:N$
 - $x(t+1) \leftarrow \text{Draw sample } x(t)$

Eric Xing

6

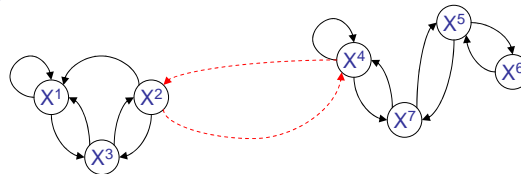
Mixing time



- The ε mixing time T_ε is the minimal number of steps (from any starting distribution) until $D_{\text{var}}(\mathcal{P}^T, \pi) \leq \varepsilon$, where D_{var} is the variational distance between the two distance:

$$D_{\text{var}}(\mu_1, \mu_2) \stackrel{\text{def}}{=} \sup_{A \subset S} |\mu_1(A) - \mu_2(A)|$$

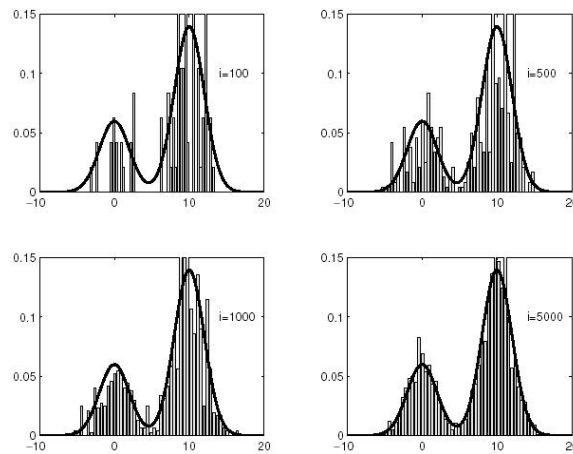
- Chains with low bandwidth (conductance) regions of space take a long time to mix.
- This arises for GMs with deterministic or highly skewed potentials.



Eric Xing

7

MCMC example



$$q(x^*|x) \sim N(x^*, 100)$$

$$p(x) \sim 0.3 \exp(-0.2x^2) + 0.7 \exp(-0.2(x-10)^2)$$

Eric Xing

8

Summary of MH



- Random walk through state space
- Can simulate multiple chains in parallel
- Much hinges on proposal distribution Q
 - Want to visit state space where $p(X)$ puts mass
 - Want $A(x^*|x)$ high in modes of $p(X)$
 - Chain mixes well
- Convergence diagnosis
 - How can we tell when burn-in is over?
 - Run multiple chains from different starting conditions, wait until they start "behaving similarly".
 - Various heuristics have been proposed.

Gibbs sampling



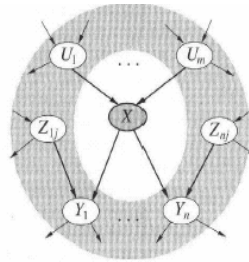
- Gibbs sampling is an MCMC algorithm that is especially appropriate for inference in graphical models.
- The procedure
 - we have variable set $\mathbf{X} = \{X_1, X_2, X_3, \dots, X_N\}$ for a GM
 - at each step one of the variables X_i is selected (at random or according to some fixed sequences), denote the remaining variables as \mathbf{X}_{-i} , and its current value as $\mathbf{x}_{-i}^{(t-1)}$
 - Using the "alarm network" as an example, say at time t we choose X_E and we denote the current value assignments of the remaining variables, \mathbf{X}_{-E} , obtained from previous samples, as $\mathbf{x}_{-E}^{(t-1)} = \{x_B^{(t-1)}, x_A^{(t-1)}, x_J^{(t-1)}, x_M^{(t-1)}\}$
 - the conditional distribution $p(X_i | \mathbf{x}_{-i}^{(t-1)})$ is computed
 - a value $x_i^{(t)}$ is sampled from this distribution
 - the sample $x_i^{(t)}$ replaces the previous sampled value of X_i in \mathbf{X} .
 - i.e., $\mathbf{x}^{(t)} = \mathbf{x}_{-E}^{(t-1)} \cup x_E^{(t)}$



Markov Blanket



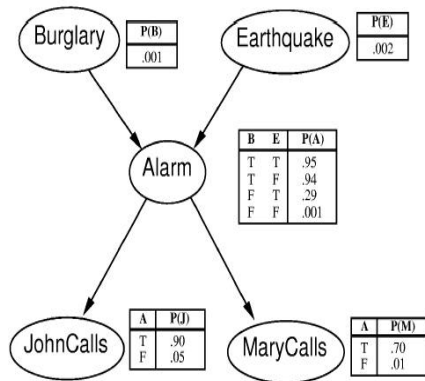
- Markov Blanket in BN
 - A variable is independent from others, given its parents, children and children's parents (d-separation).
 - MB in MRF
 - A variable is independent all its non-neighbors, given all its direct neighbors.
- $\Rightarrow P(X_i | X_{-i}) = P(X_i | MB(X_i))$
- Gibbs sampling
 - Every step, choose one variable and sample it by $P(X|MB(X))$ based on previous sample.



Eric Xing

11

Gibbs sampling of the alarm network



$MB(A) = \{B, E, J, M\}$
 $MB(E) = \{A, B\}$

- To calculate $P(J|B1, M1)$
- Choose $(B1, E0, A1, M1, J1)$ as a start
- **Evidences** are $B1, M1$, **variables** are A, E, J .
- Choose next variable as A
- Sample A by $P(A|MB(A)) = P(A|B1, E0, M1, J1)$ suppose to be false.
- $(B1, E0, A0, M1, J1)$
- Choose next random variable as E, sample $E \sim P(E|B1, A0)$
- ...

Eric Xing

12

Gibbs sampling



- Gibbs sampling is a special case of MH
- The transition matrix updates each node one at a time using the following proposal:

$$Q((x_i, \mathbf{x}_{-i}) \rightarrow (x_i', \mathbf{x}_{-i})) = p(x_i' | \mathbf{x}_{-i})$$

- This is efficient since for two reasons
 - It leads to samples that is always accepted

$$\begin{aligned} A((x_i, \mathbf{x}_{-i}) \rightarrow (x_i', \mathbf{x}_{-i})) &= \min \left(1, \frac{p(x_i' | \mathbf{x}_{-i}) Q((x_i, \mathbf{x}_{-i}) \rightarrow (x_i', \mathbf{x}_{-i}))}{p(x_i, \mathbf{x}_{-i}) Q((x_i', \mathbf{x}_{-i}) \rightarrow (x_i, \mathbf{x}_{-i}))} \right) \\ &= \min \left(1, \frac{p(x_i' | \mathbf{x}_{-i}) p(\mathbf{x}_{-i}) p(x_i | \mathbf{x}_{-i})}{p(x_i, \mathbf{x}_{-i}) p(\mathbf{x}_{-i}) p(x_i' | \mathbf{x}_{-i})} \right) = \min(1, 1) \end{aligned}$$

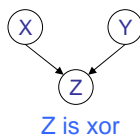
Thus $T((x_i, \mathbf{x}_{-i}) \rightarrow (x_i', \mathbf{x}_{-i})) = p(x_i' | \mathbf{x}_{-i})$

- It is efficient since $p(x_i' | \mathbf{x}_{-i})$ only depends on the values in X_i 's Markov blanket

Gibbs sampling



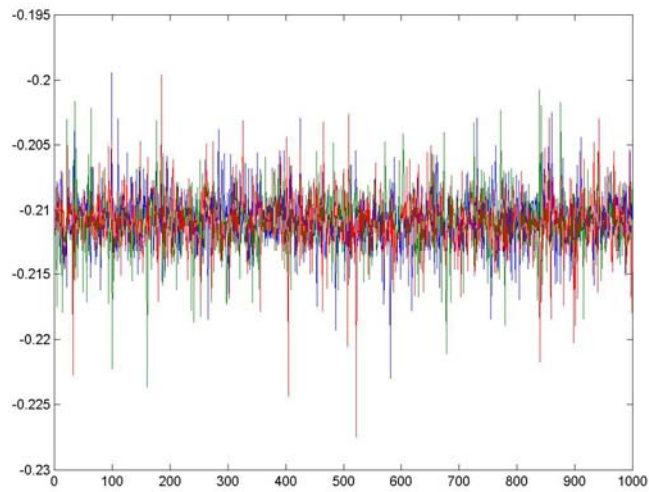
- Scheduling and ordering:
 - Sequential sweeping: in each "epoch" t , touch every r.v. in some order and yield a new sample, $\mathbf{x}^{(t)}$, after every r.v. is resampled
 - Randomly pick an r.v. at each time step
- Blocking:
 - Large state space: state vector \mathbf{X} comprised of many components (high dimension)
 - Some components can be correlated and we can sample components (i.e., subsets of r.v.s.) one at a time
- Gibbs sampling can fail if there are deterministic constraint



- Suppose we observe $Z = 1$. The posterior has 2 modes: $P(X = 1, Y = 0 | Z = 1)$ and $P(X = 0, Y = 1 | Z = 1)$. if we start in mode 1, $P(X | Y = 0, Z = 1)$ leaves $X = 1$, so we can't move to mode 2 (Reducible Markov chain).
- If all states have non-zero probability, the MC is guaranteed to be regular.
- Sampling blocks of variables at a time can help improve mixing.

GOOD!

Chains

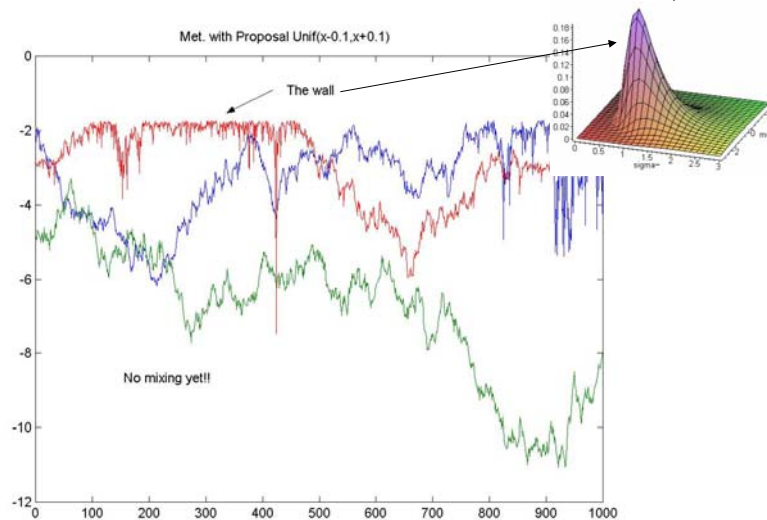


Eric Xip

15

BAD!

Chains



Eric Xip

$n=3, \alpha=1, m=0.92, s=1.55, N_{met}=5.$

16

The **Art** of simulation



- Run several chains
- Start at over-dispersed points
- Monitor the log lik.
- Monitor the serial correlations
- Monitor acceptance ratios
- Re-parameterize (to get approx. indep.)
- Re-block (Gibbs)
- Collapse (int. over other pars.)
- Run with troubled pars. fixed at reasonable vals.

Eric Xing

17

Example: Recall Latent Dirichlet Allocation

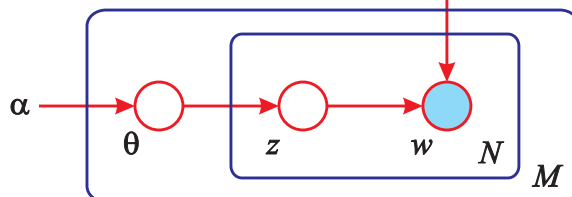


- Blei, Jordan and Ng (2003)
- Generative model of documents (but broadly applicable e.g. collaborative filtering, image retrieval, bioinformatics)
- Generative model:
 - choose θ
 - choose topic z
 - choose word w

$$\theta \sim \text{Dir}(\alpha)$$

$$z_n \sim \text{Mult}(\theta)$$

$$w_n \sim p(w_n | z_n, \beta)$$



Eric Xing

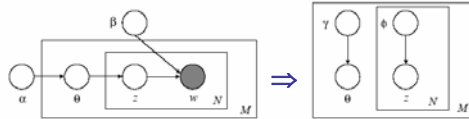
18

Variational Approximation



- Naïve Mean Field:

$$\begin{aligned}
 q(\theta, \mathbf{z}) &= q_\theta(\theta) q_z(\mathbf{z}) \\
 &= \text{Dir}(\theta \mid \gamma = f(\alpha, \langle \mathbf{z} \rangle)) \times \\
 &\quad \text{Multi}(\mathbf{z} \mid \phi = f(\beta_w, \langle \ln \theta \rangle))
 \end{aligned}$$



$$\begin{aligned}
 \phi_{ni} &\propto \beta_{i w_n} \exp\{E_q[\log(\theta_i) \mid \gamma]\} \\
 \gamma_i &= \alpha_i + \sum_{n=1}^N \phi_{ni}
 \end{aligned}$$

Eric Xing

19

Collapsed Gibbs sampling of M³ model (Tom Griffiths & Mark Steyvers)

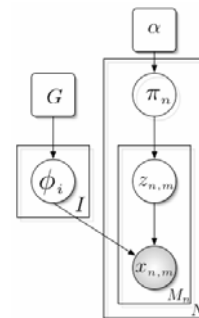


- Collapsed Gibbs sampling
 - Integrate out π

For variables $\mathbf{z} = z_1, z_2, \dots, z_n$

Draw $z_i^{(t+1)}$ from $P(z_i \mid \mathbf{z}_{-i}, \mathbf{w})$

$$\mathbf{z}_{-i} = z_1^{(t+1)}, z_2^{(t+1)}, \dots, z_{i-1}^{(t+1)}, z_{i+1}^{(t)}, \dots, z_n^{(t)}$$



Eric Xing

20

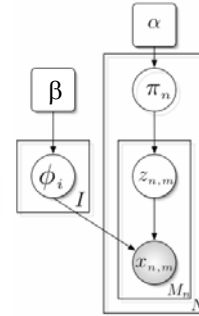
Gibbs sampling



- Need full conditional distributions for variables
- Since we only sample z we need

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto P(w_i | z_i = j, \mathbf{z}_{-i}, \mathbf{w}_{-i}) P(z_i = j | \mathbf{z}_{-i})$$

$$= \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$



$n_j^{(w)}$ number of times word w assigned to topic j
 $n_j^{(d)}$ number of times topic j used in document d

Gibbs sampling



i	w_i	d_i	iteration
			1
1	MATHEMATICS	1	2
2	KNOWLEDGE	1	2
3	RESEARCH	1	1
4	WORK	1	2
5	MATHEMATICS	1	1
6	RESEARCH	1	2
7	WORK	1	2
8	SCIENTIFIC	1	1
9	MATHEMATICS	1	2
10	WORK	1	1
11	SCIENTIFIC	2	1
12	KNOWLEDGE	2	1
.	.	.	.
.	.	.	.
50	JOY	5	2

Gibbs sampling



i	w_i	d_i	iteration	
			1	2
1	MATHEMATICS	1	z_i	z_i
2	KNOWLEDGE	1	2	?
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.
.
50	JOY	5	2	

Eric Xing

23

Gibbs sampling



i	w_i	d_i	iteration	
			1	2
1	MATHEMATICS	1	z_i	z_i
2	KNOWLEDGE	1	2	?
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.
.
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$

Eric Xing

24

Gibbs sampling



i	w_i	d_i	iteration	
			1	2
1	MATHEMATICS	1	z_i	z_i
2	KNOWLEDGE	1	2	?
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$

Eric Xing

25

Gibbs sampling



i	w_i	d_i	iteration	
			1	2
1	MATHEMATICS	1	2	2
2	KNOWLEDGE	1	2	?
3	RESEARCH	1	1	
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$

Eric Xing

26

Gibbs sampling



i	w_i	d_i	iteration	
			1	2
1	MATHEMATICS	1	2	2
2	KNOWLEDGE	1	2	1
3	RESEARCH	1	1	?
4	WORK	1	2	
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$

Eric Xing

27

Gibbs sampling



i	w_i	d_i	iteration	
			1	2
1	MATHEMATICS	1	2	2
2	KNOWLEDGE	1	2	1
3	RESEARCH	1	1	1
4	WORK	1	2	?
5	MATHEMATICS	1	1	
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$

Eric Xing

28

Gibbs sampling



<i>i</i>	w_i	d_i	iteration	
			1	2
1	MATHEMATICS	1	2	2
2	KNOWLEDGE	1	2	1
3	RESEARCH	1	1	1
4	WORK	1	2	2
5	MATHEMATICS	1	1	?
6	RESEARCH	1	2	
7	WORK	1	2	
8	SCIENTIFIC	1	1	
9	MATHEMATICS	1	2	
10	WORK	1	1	
11	SCIENTIFIC	2	1	
12	KNOWLEDGE	2	1	
.	.	.	.	
.	.	.	.	
50	JOY	5	2	

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$

Eric Xing

29

Gibbs sampling



<i>i</i>	w_i	d_i	iteration		...	1000
			1	2		
1	MATHEMATICS	1	2	2		z_i
2	KNOWLEDGE	1	2	1		2
3	RESEARCH	1	1	1		2
4	WORK	1	2	2		1
5	MATHEMATICS	1	1	2		2
6	RESEARCH	1	2	2		2
7	WORK	1	2	2		2
8	SCIENTIFIC	1	1	1	...	1
9	MATHEMATICS	1	2	2		2
10	WORK	1	1	2		2
11	SCIENTIFIC	2	1	1		2
12	KNOWLEDGE	2	1	2		2
.
.
50	JOY	5	2	1		1

$$P(z_i = j | \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T\alpha}$$

Eric Xing

30

Document tagging



"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.