

Reasoning (IV)

Mingsheng Long

Tsinghua University

Outline

- Sampling Method

- Monte Carlo

- Markov Chain Monte Carlo

- Metropolis-Hastings Algorithm

- Gibbs Sampling

- Sampling Method for LDA

- Bayesian Estimation

- Dirichlet-Multinomial Conjugate

- Gibbs Sampling for LDA



ELBO and EM Algorithm

- Evidence Lower Bound (ELBO):

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_z q(z) \log \left(\frac{p(x, z | \theta)}{q(z)} \right) \\ &= -\text{KL}[q(z) || p(z|x, \theta)] + \log p(x|\theta)\end{aligned}$$

- EM Algorithm:

- Choose initial θ^{old}

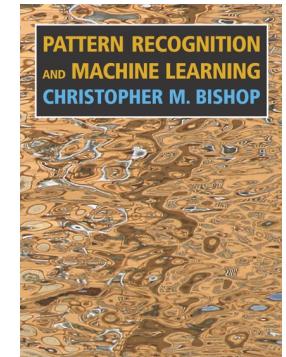
- E-step: Let $q^* = \underset{q}{\operatorname{argmax}} \mathcal{L}(q, \theta^{\text{old}})$

Minimize

Optimal choice for E-step
 $q^*(z) = p(z|x, \theta^{\text{old}})$

- M-step: Let $\theta^{\text{new}} = \underset{\theta}{\operatorname{argmax}} \mathcal{L}(q^*, \theta)$

- Go to step 2, until converged.



PRML

Chapter 11

Intractable E-Step

- We need to solve: Optimal choice for E-step $q^*(z) = p(z|x, \theta^{\text{old}})$

$$\text{In LDA: } p(\theta, z, \beta | w, \alpha, \eta) = \frac{p(\theta, z, \beta, w | \alpha, \eta)}{p(w | \alpha, \eta)}$$

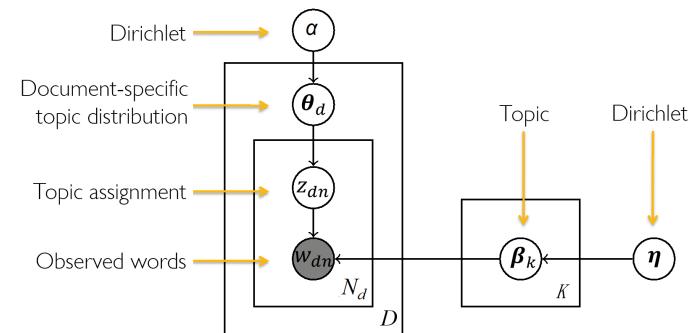
- But the denominator is **intractable!**

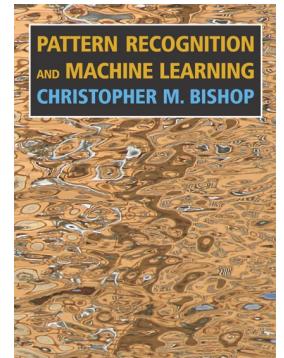
$$p(w|\alpha, \eta) = \int \int \sum_z p(\theta, z, \beta, w | \alpha, \eta) d\theta d\beta$$

- We need **approximation** here!

- There are two ways to solve E-step approximately:

- Variational Inference (last lecture)
 - Sampling – **Markov Chain Monte Carlo** (this lecture)





Sampling Method

- We want to estimate the integral on some distribution:

$$\int_{x \in \mathcal{X}} f(x)p(x)dx = \mathbb{E}_p f(x)$$

- Sampling method:

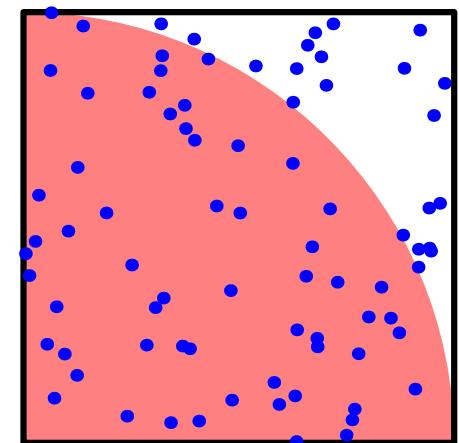
PRML
Chapter 11

- Draw an i.i.d. sample set from distribution p : $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$
- In practice: $\{\mathbf{x}_{1:n}\}$ not independent, effective sample size is small.
- We can only generate uniform pseudo-random numbers in $U[0,1]$.
- Monte Carlo Estimator: $\frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i)$
 - Unbiased: $\mathbb{E}_{\mathcal{D} \sim p^n} \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbf{x}_i \sim p} f(\mathbf{x}_i) = \mathbb{E}_p f(x)$
 - Variance: $\text{Var}_{\mathcal{D} \sim p^n} \left[\frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i) \right] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}_p [f(\mathbf{x}_i)] = \frac{1}{n} \text{Var}_p [f]$

Monte Carlo Method

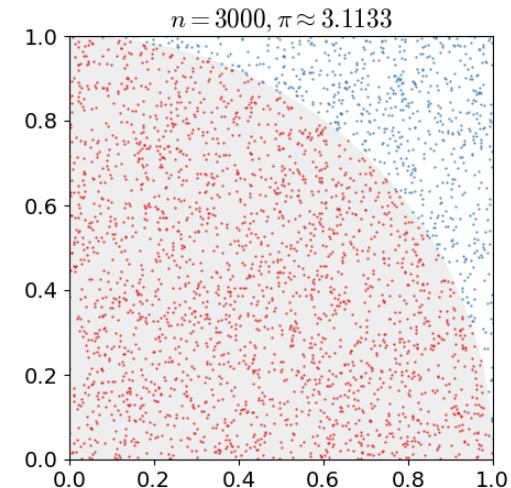
- Many problems can be changed into integral form.
- Example: the computation of π .
 - View π as a quarter of area of a unit circle.

$$\pi = 4 \iint \mathbb{I}(x^2 + y^2 < 1) P(x, y) dx dy$$



- Consider a uniform distribution on $U[0,1]^2$

$$P(x, y) = \begin{cases} 1 & 0 < x < 1 \text{ and } 0 < y < 1 \\ 0 & \text{otherwise} \end{cases}$$



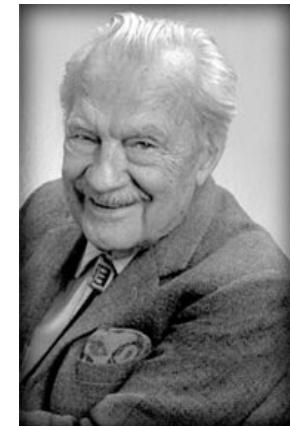
Monte Carlo Method



Stanislaw
Ulam



John von
Neumann



Nicholas
Metropolis

- For being secret, this work required a **code name**.
 - Monte Carlo is a city in Monaco.
- “Monte Carlo is an **extremely bad method**; it should be used only when **all alternative methods are worse**.” — Alan Sokal, 1996



Monte Carlo EM Algorithm

- Evidence lower-bound (ELBO):

$$\mathcal{L}(q, \theta) = \sum_z q(z) \log \left(\frac{p(x, z | \theta)}{q(z)} \right) = \mathbb{E}_{z \sim q} \log \left(\frac{p(x, z | \theta)}{q(z)} \right)$$

- EM Algorithm:

- E-step: $q^* = \operatorname{argmax}_q \mathcal{L}(q, \theta^{\text{old}}) = p(z|x, \theta^{\text{old}})$

Sampling from it!

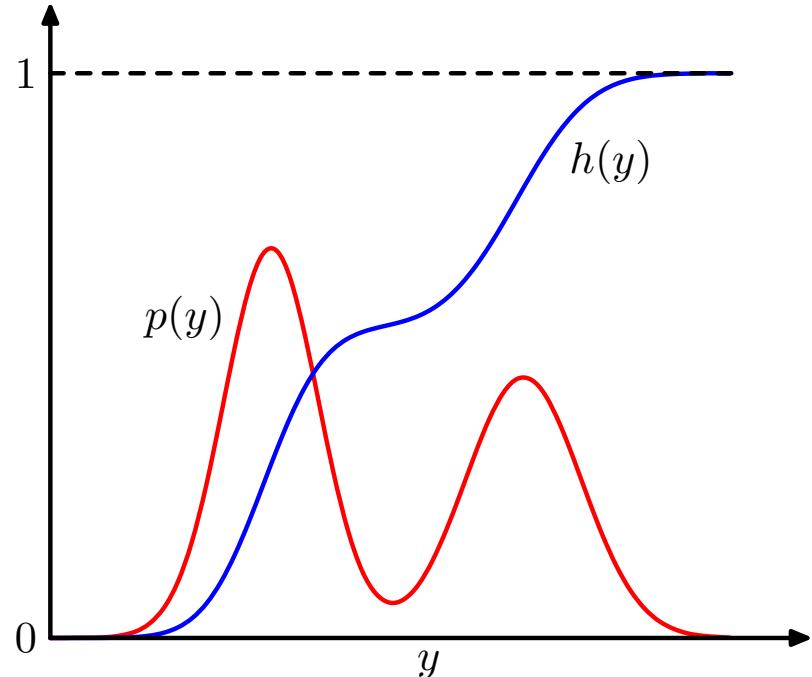
- M-Step: $\theta^{\text{new}} = \operatorname{argmax}_{\theta} \mathcal{L}(q^*, \theta) = \operatorname{argmax}_{\theta} \mathbb{E}_{z \sim q} \log \left(\frac{p(x, z | \theta)}{q(z)} \right)$

$$= \operatorname{argmax}_{\theta} \mathbb{E}_{z \sim q} \log p(x, z | \theta) \approx \operatorname{argmax}_{\theta} \frac{1}{T} \sum_{i=1}^T \log p(x, z_i | \theta)$$

- We don't need to explicitly compute q^* but sample z_1, \dots, z_T from q^* !



Sampling from Distribution

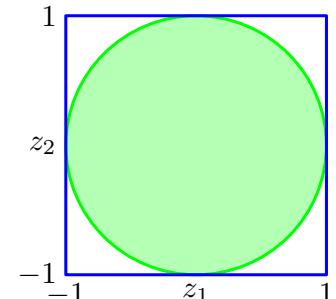
- Key problem: How to sample from some complex distribution p ?
 - Suppose we know the sampling method of some simple distribution.
 - We convert samples to target distributions p from the simple one.
 - For example, if we can sample from
 - Uniform distribution on $[0,1]$.
 - For a target p with explicit CDF:
 - $h(x) = \int_{-\infty}^x p(z) dz$
- First, sample $u \sim U[0,1]$.
 - Second, compute $x(u) = h^{-1}(u)$
- 

Example: Sampling from Gaussian

- Box-Muller method: generate samples from Gaussian distribution
 - Generate pairs of uniform random numbers $z_1, z_2 \in (-1,1)$
 - Discard each pair unless it satisfies $z_1^2 + z_2^2 \leq 1$
 - A uniform distribution inside the unit circle with $p(z_1, z_2) = \frac{1}{\pi}$

- Define $r^2 = z_1^2 + z_2^2$, and

$$y_1 = z_1 \left(\frac{-2 \ln r^2}{r^2} \right)^{1/2}, \quad y_2 = z_2 \left(\frac{-2 \ln r^2}{r^2} \right)^{1/2}$$



- Then
$$p(y_1, y_2) = p(z_1, z_2) \left| \frac{\partial(z_1, z_2)}{\partial(y_1, y_2)} \right|$$
$$= \left[\frac{1}{\sqrt{2\pi}} \exp(-y_1^2/2) \right] \left[\frac{1}{\sqrt{2\pi}} \exp(-y_2^2/2) \right]$$

Rejection Sampling

When unable to **compute**
and **invert** $h(x)$.



Sampling based on $p(x)$.
Easy to evaluate $p(x) = \frac{\tilde{p}(x)}{Z}$

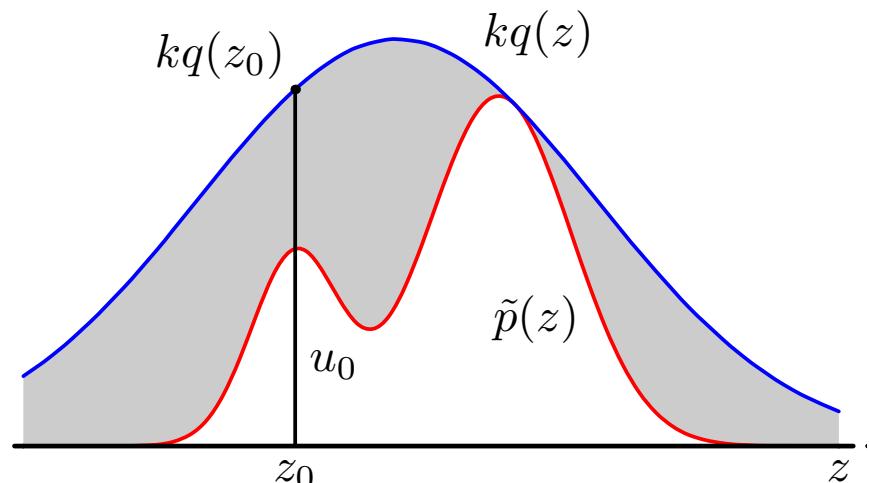
unknown

- Find a proposal distribution q :
 - We can sample from it.
 - $\exists k, k \cdot q(x) \geq p(x)$

- Rejection Sampling:
 - Draw $x \sim q, u \sim U[0, k \cdot q(x)]$

- - Reject x if $u > p(x)$.

- $p(x|\text{accept}) \propto q(x)p(\text{accept}|x) = q(x) \left(\frac{p(x)}{k \cdot q(x)} \right) = \frac{p(x)}{k}$
- The acceptance rate decreases **exponentially** with dimensionality of x .



Outline

- **Sampling Method**
 - Monte Carlo
 - **Markov Chain Monte Carlo**
 - Metropolis-Hastings Algorithm
 - Gibbs Sampling
- Sampling Method for LDA
 - Bayesian Estimation
 - Dirichlet-Multinomial Conjugate
 - Gibbs Sampling for LDA



Monte Carlo for LDA

- We want to compute the following terms in E-step of LDA:

$$p(z | w, \theta, \beta) = \frac{p(z, w | \beta, \theta)}{p(w | \beta, \theta)}$$

Easier to
compute!

Intractable!
No explicit PDF.

- The problem turns into a canonical problem:

- If we want to sample from a distribution:

$$p(x = j) = \pi_j \propto b_j$$

- Where b_j is computable, the value of $\sum_{j=1}^m b_j = B$ is unknown.

- Question: How to sample from distribution $p(x)$ for Monte Carlo?



Markov Chain

- Markov Chain:

$$P(x_{t+1}|x_1, \dots, x_t) = P(x_{t+1}|x_t)$$

- Assume that the state space is finite:

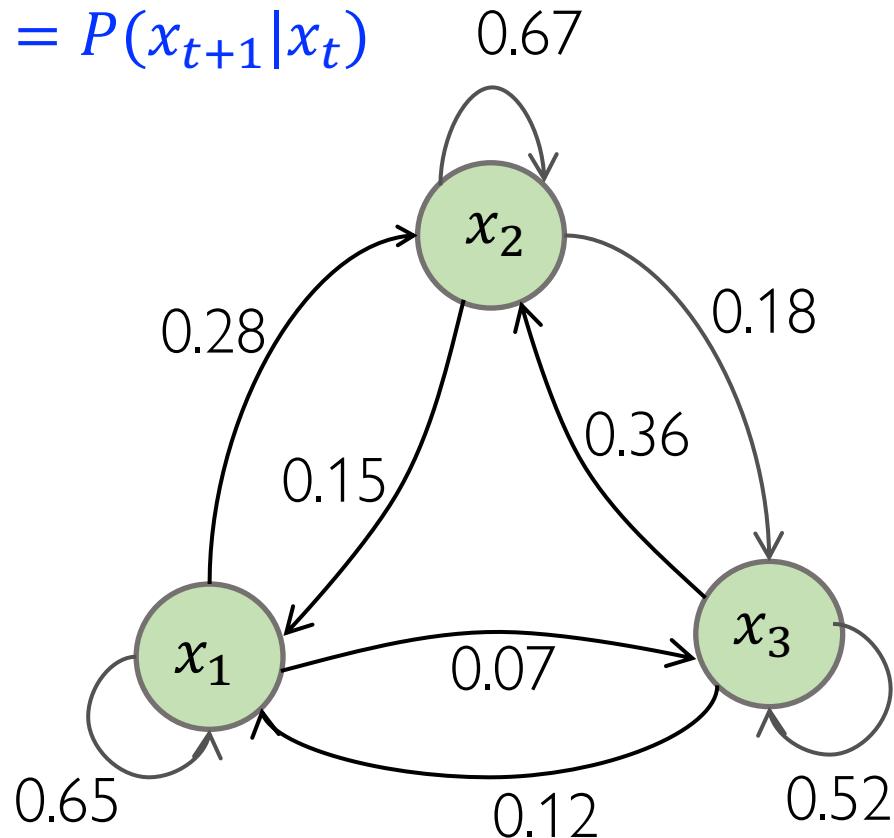
- $\mathcal{X} = \{1, \dots, k\}$

- Transition matrix P :

- $P_{ij} = P(x_{t+1} = j | x_t = i)$

- For example:

$$P = \begin{bmatrix} 0.65 & 0.28 & 0.07 \\ 0.15 & 0.67 & 0.18 \\ 0.12 & 0.36 & 0.52 \end{bmatrix}$$



- Describe the probability of transition from one state to other states.

Convergence of Markov Chain

- We use **state distribution** to represent probability of state at time t :

$$\boldsymbol{\pi}^t = (\pi_1^t \quad \dots \quad \pi_k^t)$$

- One transition: $\boldsymbol{\pi}^{t+1} = \boldsymbol{\pi}^t \mathbf{P}$

$$p(x_{t+1} = i) = \sum_{j=1}^k p(x_t = j) P_{ji}$$

- An example:

- If $\boldsymbol{\pi}^0 = (0.21 \quad 0.68 \quad 0.11)$ and transit for T times:

$$\boldsymbol{\pi}^T = \boldsymbol{\pi}^0 (\mathbf{P})^T$$

- The result will be $\boldsymbol{\pi}^{10} = (0.286 \quad 0.489 \quad 0.225)$

- If we change the **initial state**, the result will **remain unchanged**.

- Observation: $\boldsymbol{\pi}^{10} \mathbf{P} = \boldsymbol{\pi}^{10}$.



Convergence of Markov Chain

- **Ergodicity:** A Markov chain has a unique stationary distribution if:
 - From all initial states π^0 , the chain will converge to π :
$$\pi = \pi P$$
 - After a sufficient number of transitions.
 - Each transition is equivalent to sampling from distribution π .

- **Detailed balance condition**
 - Markov chain is ergodic if and only if: $\pi_i P_{ij} = \pi_j P_{ji}$

$$\sum_{i=1}^{+\infty} \pi_i P_{ij} = \sum_{i=1}^{+\infty} \pi_j P_{ji} = \pi_j \sum_{i=1}^{+\infty} P_{ji} = \pi_j$$



Sampling with Markov Chain

- The problem turns into a canonical problem:

- If we want to sample from a distribution:

$$p(x = j) = \pi_j \propto b_j$$

- Where b_j is computable, the value of $\sum_{j=1}^m b_j = B$ is unknown.

- Question: How to sample from distribution $p(x)$ for Monte Carlo?

- Answer: Find some \mathbf{P} such that the stationary distribution is $\pi_j \propto b_j$:

- We use random samples transiting with \mathbf{P} :
 - After several transitions, the state distribution will converge to π .
 - The transiting samples after convergence will follow distribution π .



Markov Chain Monte Carlo

- Given a Markov chain with known transition matrix \mathbf{P} (unknown!).
 - Assume that detailed balance condition is hold: $\pi_i P_{ij} = \pi_j P_{ji}$

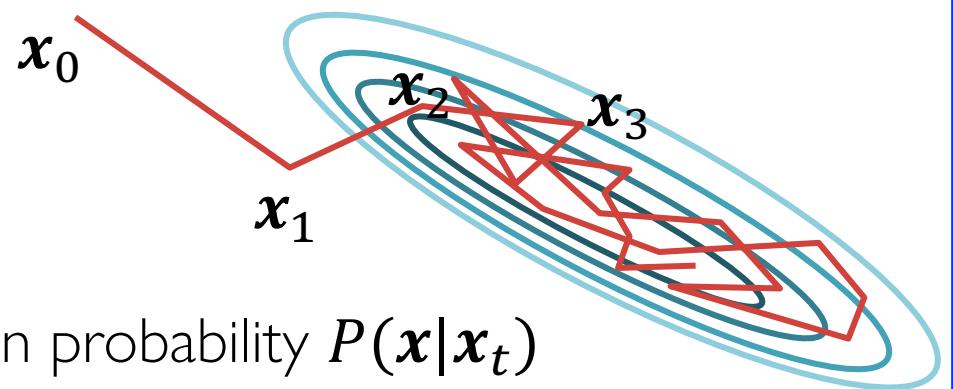
- Markov chain sampling:

- Start from random \mathbf{x}_0

- For $t = 0, \dots, T - 1$

- Sample \mathbf{x}_{t+1} from transition probability $P(\mathbf{x}|\mathbf{x}_t)$

- Assume convergence at time T_0



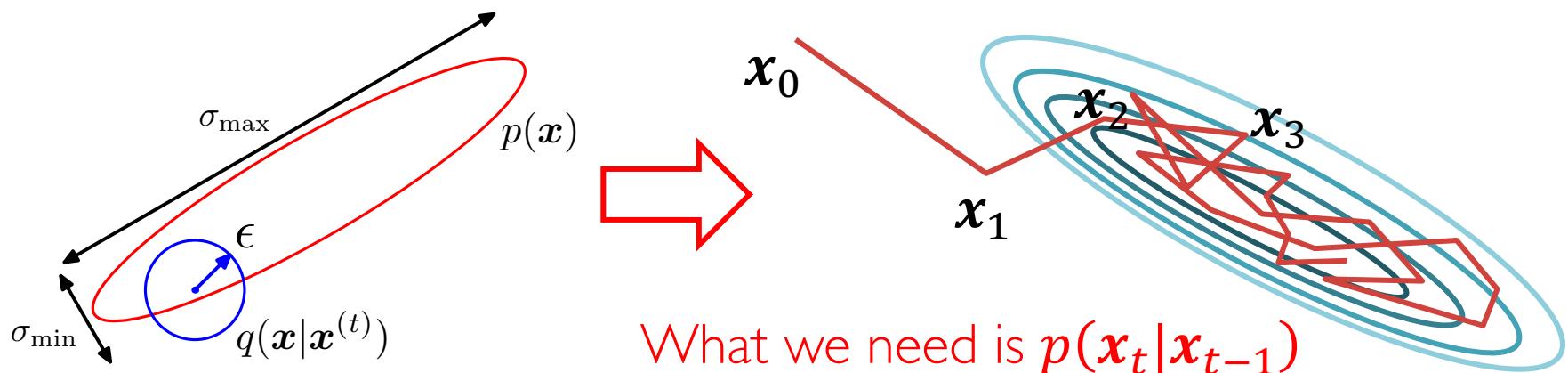
- Then the sample set $\{\mathbf{x}_{T_0}, \dots, \mathbf{x}_T\}$ follows stationary distribution $\boldsymbol{\pi}$.

- Monte Carlo with Markov Chain:

$$\frac{1}{T-T_0} \sum_{t=T_0}^{T-1} f(\mathbf{x}_t) \rightarrow \mathbb{E}_{\mathbf{x} \sim \boldsymbol{\pi}} f(\mathbf{x})$$

Markov Chain Monte Carlo

- We do not know transition matrix \mathbf{P} in practice. How to sample?
- We first choose a proposal distribution q as transition probability:
 - Usually set as Gaussian: $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t | \mathbf{x}_{t-1}, \epsilon^2)$.
 - This distribution is symmetric: $q(\mathbf{x}|\mathbf{y}) = q(\mathbf{y}|\mathbf{x})$
 - Forms Markov Chain with transition matrix \mathbf{Q} : $Q_{ij} = \mathcal{N}(\mathbf{x}_j | \mathbf{x}_i, \epsilon^2)$



Markov Chain Monte Carlo

- Transition matrix by proposal distribution \mathbf{Q} : $Q_{ij} = \mathcal{N}(x_j | x_i, \epsilon^2)$

- Detailed balance condition may be violated:

$$\pi_i Q_{ij} \neq \pi_j Q_{ji}$$

- We add acceptance rate α_{ij} and α_{ji} to keep it balance:

$$\pi_i Q_{ij} \alpha_{ij} = \pi_j Q_{ji} \alpha_{ji}$$

- Let new transition matrix $P_{ij} = Q_{ij} \alpha_{ij}$ and $P_{ji} = Q_{ji} \alpha_{ji}$

- Then we have $\pi_i P_{ij} = \pi_j P_{ji}$.

- Detailed balance condition holds! We can use it to sample.

- Given transition matrix \mathbf{Q} , how to make Markov chain transit with \mathbf{P} ?

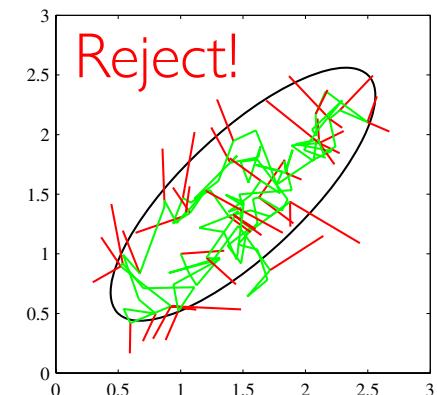


Markov Chain Monte Carlo

- Given transition matrix \mathbf{Q} , how to make Markov chain transit with \mathbf{P} ?

- Rejection Sampling!

- Sample $\mathbf{x}_* \sim q(\mathbf{x}_t | \mathbf{x}_{t-1})$.
 - Suppose $\mathbf{x}_{t-1} = \mathbf{x}_i, \mathbf{x}_* = \mathbf{x}_j$.
- Sample $u \sim U[0,1]$.
- If $u > \alpha_{ij}$, **reject** the transition: $\mathbf{x}_t = \mathbf{x}_{t-1}$.
- Else, **accept**: $\mathbf{x}_t = \mathbf{x}_*$.



- Assume convergence at time T_0
- Then the sample set $\{\mathbf{x}_{T_0}, \dots, \mathbf{x}_T\}$ follows stationary distribution $\boldsymbol{\pi}$.

Outline

- **Sampling Method**

- Monte Carlo

- Markov Chain Monte Carlo

- **Metropolis-Hastings Algorithm**

- Gibbs Sampling

- Sampling Method for LDA

- Bayesian Estimation

- Dirichlet-Multinomial Conjugate

- Gibbs Sampling for LDA



Metropolis-Hastings Algorithm

- Acceptance rate α_{ij} : the larger, the more efficient.
- Recall the detailed balance condition: $\pi_i Q_{ij} \alpha_{ij} = \pi_j Q_{ji} \alpha_{ji}$
- A natural choice: $\alpha_{ij} = \pi_j Q_{ji}$ and $\alpha_{ji} = \pi_i Q_{ij}$
 - Then we have: $\pi_i Q_{ij} \alpha_{ij} = \pi_i Q_{ij} \pi_j Q_{ji} = \pi_j Q_{ji} \alpha_{ji}$.
- Two problems:
 - First, π_j is not computable. Second, $\pi_j Q_{ji}$ is too small, very slow.
- Final choice: $\alpha_{ij} = \min(1, \frac{\pi_j Q_{ji}}{\pi_i Q_{ij}})$. Big enough to accept more.
 - Compute α_{ij} only needs $\frac{\pi_j}{\pi_i} = \frac{b_j}{b_i}$. No need of intractable B .
 - $\pi_i Q_{ij} \alpha_{ij} = \min(\pi_i Q_{ij}, \pi_j Q_{ji}) = \pi_j Q_{ji} \alpha_{ji}$. Detailed balance holds.



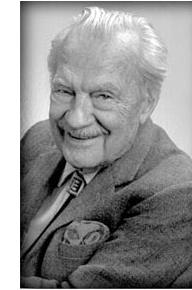
Metropolis-Hastings Algorithm

- Metropolis-Hastings Algorithm

- Sample $\mathbf{x}_* \sim q(\mathbf{x}_t | \mathbf{x}_{t-1})$.

- Suppose $\mathbf{x}_{t-1} = \mathbf{x}_i, \mathbf{x}_* = \mathbf{x}_j$.

- Sample $u \sim U[0,1]$.



Nicholas
Metropolis



Wilfred
Keith Hastings

- If $u > \alpha_{ij} = \min\left(1, \frac{\pi_j Q_{ji}}{\pi_i Q_{ij}}\right)$, **reject** the transition: $\mathbf{x}_t = \mathbf{x}_{t-1}$.
- Else, **accept**: $\mathbf{x}_t = \mathbf{x}_*$.

- Solves the problem of sampling from **any** distribution for Monte Carlo!
- Chosen as one of the top 10 algorithms in the 20th century.

Outline

- **Sampling Method**
 - Monte Carlo
 - Markov Chain Monte Carlo
 - Metropolis-Hastings Algorithm
 - **Gibbs Sampling**
- Sampling Method for LDA
 - Bayesian Estimation
 - Dirichlet-Multinomial Conjugate
 - Gibbs Sampling for LDA



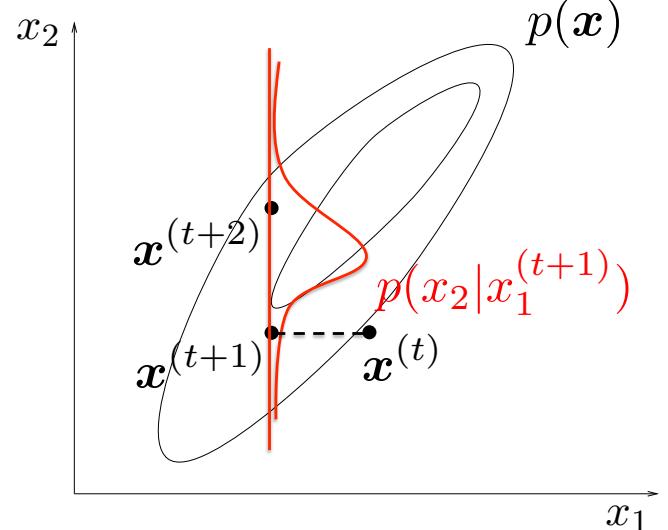
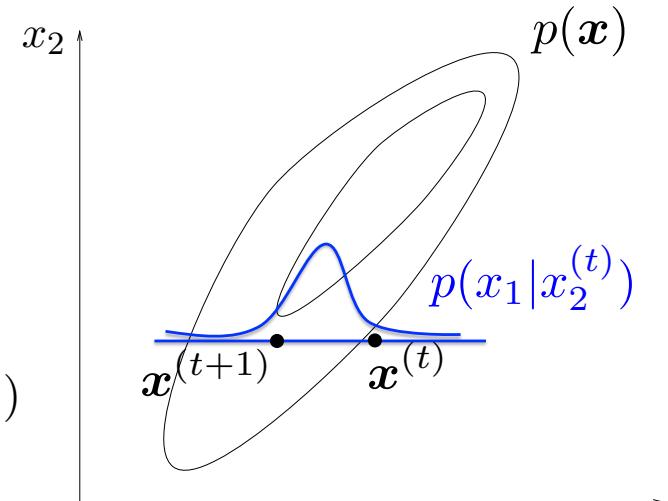
Gibbs Sampling

- Problems with Metropolis-Hastings Algorithm
 - Computing $\alpha_{ij} = \min\left(1, \frac{\pi_j Q_{ji}}{\pi_i Q_{ij}}\right)$ is inefficient in **high-dimensions**
 - $\alpha_{ij} \leq 1$, so rejection rate may high.
 - Can we **not reject?**
 - Sometimes it is impossible to compute the **joint distribution** in high-dimensions
 - Can we sample only based on the **conditional distributions** between different dimensions?
- The answer: **Gibbs sampling**. The algorithm for the big data era!



Gibbs Sampling

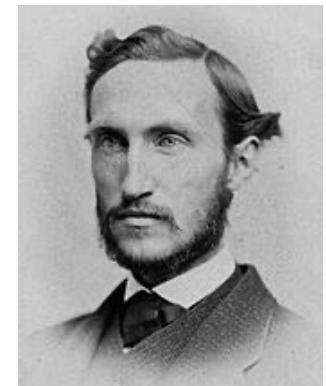
- 1. $x_1^{(t+1)} \sim p(x_1|x_2^{(t)}, x_3^{(t)}, \dots, x_n^{(t)})$
- 2. $x_2^{(t+1)} \sim p(x_2|x_1^{(t+1)}, x_3^{(t)}, \dots, x_n^{(t)})$
-
- 3. ... Sample Fix the rest $n - 1$ var.
- 4. $x_j^{(t+1)} \sim p(x_j|x_1^{(t+1)}, \dots, x_{j-1}^{(t+1)}, x_{j+1}^{(t)}, \dots, x_n^{(t)})$
- 5. ...
- 6. $x_n^{(t+1)} \sim p(x_n|x_1^{(t+1)}, x_2^{(t)}, \dots, x_{n-1}^{(t+1)})$



- Iterate across dimensions until convergence.
- One of the most beautiful algorithms in 20th century.

Gibbs Sampling

- Gibbs Sampling is a **special case** of Metropolis-Hastings Algorithm
 - The **proposal distributions** are **posterior conditionals** (computable!)
 - Accept all samples $\alpha_{ij} = 1$
- Proof:
 - We sample from the **conditional** $p(x_k^* | \mathbf{x}_{\neg k})$
 - Note that in each iteration we fix $\mathbf{x}_{\neg k}^* = \mathbf{x}_{\neg k}$
 - Transition probability $q(\mathbf{x}^* | \mathbf{x}) = p(x_k^* | \mathbf{x}_{\neg k})$
 - Joint distribution $p(\mathbf{x}) = p(\mathbf{x}_{\neg k})p(x_k | \mathbf{x}_{\neg k})$
 - Accept: $\alpha_k = \frac{\pi_j Q_{ji}}{\pi_i Q_{ij}} = \frac{p(\mathbf{x}^*) q(\mathbf{x} | \mathbf{x}^*)}{p(\mathbf{x}) q(\mathbf{x}^* | \mathbf{x})} = \frac{p(\mathbf{x}_{\neg k}^*) p(x_k^* | \mathbf{x}_{\neg k}^*) p(x_k | \mathbf{x}_{\neg k}^*)}{p(\mathbf{x}_{\neg k}) p(x_k | \mathbf{x}_{\neg k}) p(x_k^* | \mathbf{x}_{\neg k})} = 1$



Josiah W. Gibbs

Outline

- Sampling Method
 - Monte Carlo
 - Markov Chain Monte Carlo
 - Metropolis-Hastings Algorithm
 - Gibbs Sampling
- Sampling Method for LDA
 - Bayesian Estimation
 - Dirichlet-Multinomial Conjugate
 - Gibbs Sampling for LDA



Bayes Rule

- Bayes rule:

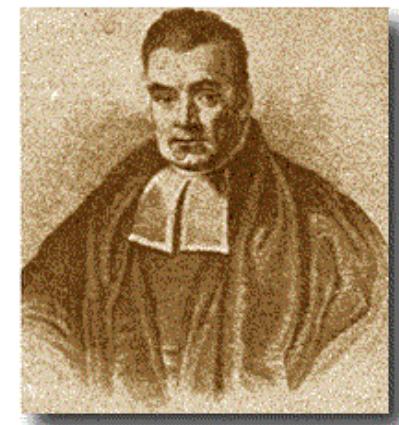
Observed Data Parameter

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

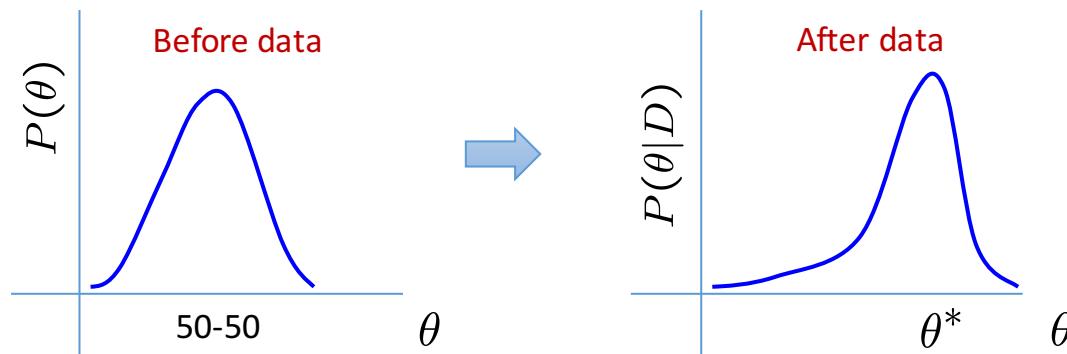
- Or equivalently:

$$p(\theta|D) \propto p(D|\theta)p(\theta)$$

↓ ↓ ↓
Posterior Likelihood Prior



Thomas Bayes

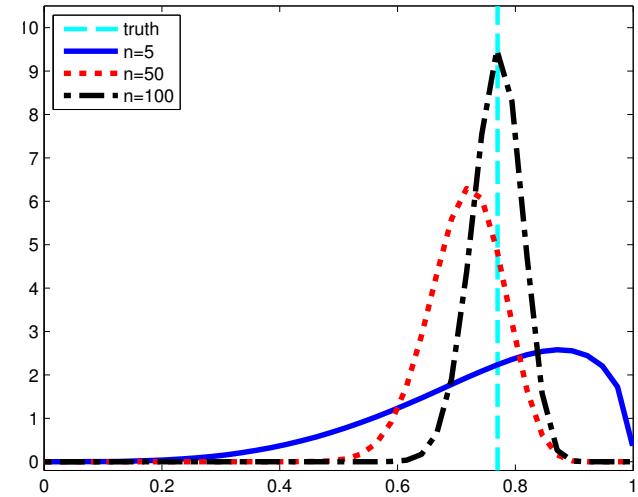


Maximum A Posteriori Estimation

$$\hat{\theta} = \arg \max_{\theta} \log p(\theta|D) = \arg \max_{\theta} \{ \log p(D|\theta) + \log p(\theta) \}$$

↓
Posterior ↓ Likelihood ↓
Prior

- MAP: Maximum a posteriori estimation of parameters θ
- We can view MLE as MAP
 - with a uniform prior distribution.
- As amount of data becomes large, posterior variance becomes small, and MAP behaves like other estimators such as MLE.

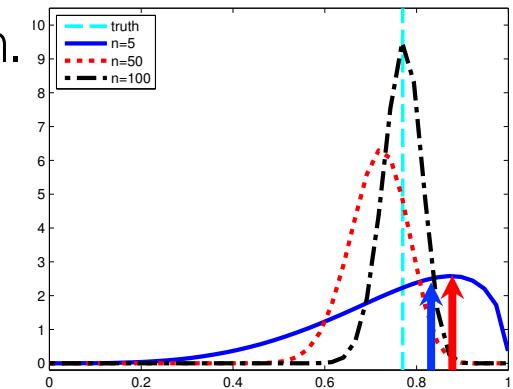


Bayes Estimator

$$\hat{\theta} = \mathbb{E}_{\theta \sim p(\theta|D)}[\theta] = \int \theta \cdot p(\theta|D)d\theta$$

↓ ↓
Mean Posterior

- Bayes Estimator of parameters θ
 - Posterior $p(\theta|D)$ is a distribution of parameter.
 - The mean of this distribution is a stable solution.
- The solutions to MAP and Bayes Estimator are not distant when distribution is unimodal.
 - Choose one by mathematical convenience.



Outline

- Sampling Method

- Monte Carlo

- Markov Chain Monte Carlo

- Metropolis-Hastings Algorithm

- Gibbs Sampling

- Sampling Method for LDA

- Bayesian Estimation

- **Dirichlet-Multinomial Conjugate**

- Gibbs Sampling for LDA



Multinomial Distribution

- Multinomial Distribution models the probability of counts of each side for rolling a k -sided dice n times.

- Each side equips with a probability θ_i .

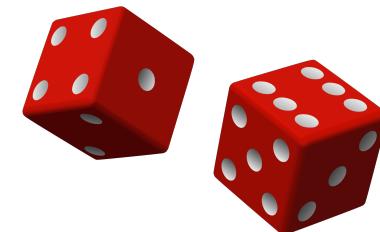
- $\sum_{i=1}^k \theta_i = 1$.

- Write by $\text{Mult}(n, \boldsymbol{\theta})$

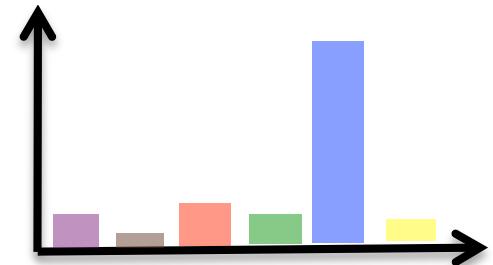
- $p(n_1, n_2, \dots, n_k) = \frac{n!}{n_1! \dots n_k!} \theta_1^{n_1} \theta_2^{n_2} \dots \theta_k^{n_k}$

- $\sum_{i=1}^k n_i = n$

- When $n = 1$, it is Categorical Distribution,
 $\text{Cat}(\boldsymbol{\theta}) = \text{Mult}(1, \boldsymbol{\theta})$



Flipping dice:
 $\sum_{l=1}^6 \theta_l = 1$

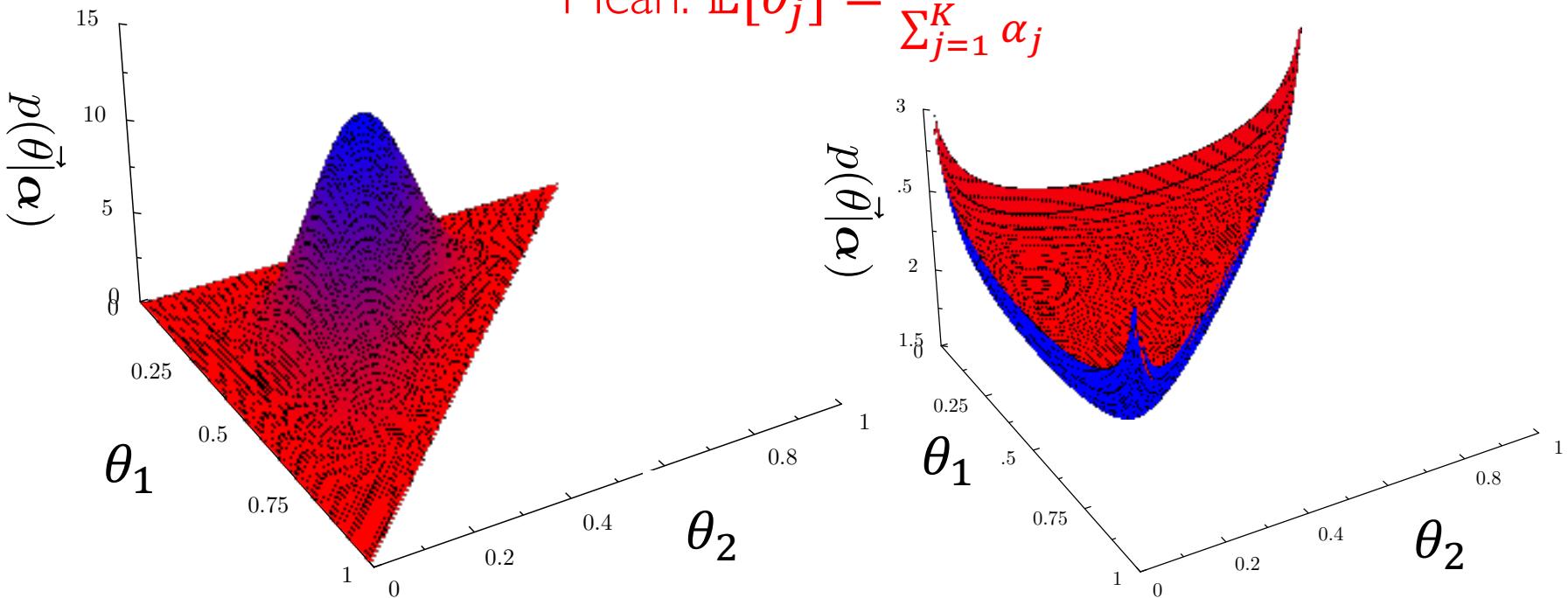


Dirichlet-Multinomial Model

- Dirichlet Distribution: Multi-dimensional version of Beta Distribution

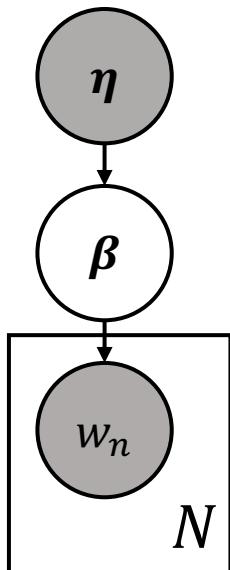
$$p(\vec{\theta}|\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k-1} \quad \text{Where } B(\alpha) = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)}$$

Mean: $\mathbb{E}[\theta_j] = \frac{\alpha_j}{\sum_{j=1}^K \alpha_j}$



Dirichlet-Multinomial Model

- Consider a simple document generation setting.
 - $\beta \sim \text{Dir}(\beta|\eta)$, dimension of β is V and $w \sim \text{Mult}(1, \beta)$.
- We can observe N words.
 - The count of term i is n_i . ($1 \leq i \leq V, \sum_{i=1}^V n_i = N$)



the	he	she	
w_1	w_2	w_3	
the	and	the	
w_4	w_5	w_6	
the	she	she	is
w_7	w_8	w_9	w_{10}

No.	Term	Count
1	the	4
2	he	1
3	she	3
4	and	1
5	...	

Posterior Distribution

- Compute the posterior distribution $p(\beta|\eta, \mathcal{W})$
 - Where \mathcal{W} is the whole corpus.
- So that is to compute:

$$p(\beta|\eta, \mathcal{W}) = \frac{p(\mathcal{W}|\beta, \eta)p(\beta|\eta)}{p(\mathcal{W}|\eta)} = \frac{p(\mathcal{W}|\beta)p(\beta|\eta)}{p(\mathcal{W}|\eta)}$$

$$\propto p(\mathcal{W}|\beta)p(\beta|\eta) = \prod_{i=1}^N p(w_i|\beta) p(\beta|\eta) = \prod_{i=1}^N \beta_{w_i} \frac{1}{B(\eta)} \prod_{j=1}^V \beta_j^{\eta_j - 1}$$

$$\prod_{i=1}^N \beta_{w_i} = \prod_{j=1}^V \beta_j^{\sum_{i=1}^N \mathbf{1}\{w_i=j\}} = \prod_{j=1}^V \beta_j^{n_j} \quad n_j: \text{Count of } j^{th} \text{ term}$$



Posterior Distribution

- Finally, we have:

$$\begin{aligned} p(\boldsymbol{\beta}|\boldsymbol{\eta}, \mathcal{W}) &\propto \prod_{i=1}^N \beta_{w_i} \frac{1}{B(\boldsymbol{\eta})} \prod_{j=1}^V \beta_j^{\eta_j-1} = \prod_{j=1}^V \beta_j^{n_j} \frac{1}{B(\boldsymbol{\eta})} \prod_{j=1}^V \beta_j^{\eta_j-1} \\ &\propto \prod_{j=1}^V \beta_j^{\eta_j+n_j-1} \end{aligned}$$

- Recall that $\text{Dir}(\boldsymbol{\beta}|\boldsymbol{\eta}) = \frac{1}{B(\boldsymbol{\eta})} \prod_{j=1}^V \beta_j^{\eta_j-1} \propto \prod_{j=1}^V \beta_j^{\eta_j-1}$

- We find that $p(\boldsymbol{\beta}|\boldsymbol{\eta}, \mathcal{W})$ is also a Dirichlet Distribution:

$$p(\boldsymbol{\beta}|\boldsymbol{\eta}, \mathcal{W}) = \text{Dir}(\boldsymbol{\beta}|\boldsymbol{\eta} + \mathbf{n})$$

Conjugate!



Dirichlet-Multinomial Model

- The Dirichlet is **conjugate** to the Multinomial:

$$\boldsymbol{\beta} \sim \text{Dir}(\boldsymbol{\eta})$$

[draw distribution over words]

For each word $n \in \{1, \dots, N\}$

$$w_n \sim \text{Mult}(1, \boldsymbol{\beta})$$

[draw word]

- The prior distribution of $\boldsymbol{\beta}$ is a **Dirichlet distribution** $\text{Dir}(\boldsymbol{\beta} | \boldsymbol{\eta})$.
- Observe some samples from **Multinomial** parameterized by $\boldsymbol{\beta}$.
- Then the posterior is also a **Dirichlet distribution**:
$$p(\boldsymbol{\beta} | \boldsymbol{\eta}, \mathcal{W}) \sim \text{Dir}(\boldsymbol{\beta} | \boldsymbol{\eta} + \mathbf{n})$$
 - The difference depends on **term counts** \mathbf{n} .
 - Will see **conjugate** distribution brings great convenience to inference.

Conjugate!



MAP vs. Bayes Estimator

- Now we get our estimation $p(\boldsymbol{\beta}|\boldsymbol{\eta}, \mathcal{W}) \sim \text{Dir}(\boldsymbol{\beta}|\boldsymbol{\eta} + \mathbf{n})$.
- Maximum a Posterior:

$$\widehat{\boldsymbol{\beta}}_{\text{MAP}} = \operatorname{argmax}_{\boldsymbol{\beta}} \text{Dir}(\boldsymbol{\beta}|\boldsymbol{\eta} + \mathbf{n})$$

$$(\widehat{\boldsymbol{\beta}}_{\text{MAP}})_i = \frac{\eta_i + n_i - 1}{\sum_{j=1}^V (\eta_j + n_j - 1)} = \frac{\eta_i + n_i - 1}{\sum_{j=1}^V \eta_j + N - V}$$

- Bayes Estimator:

May < 0

$$\widehat{\boldsymbol{\beta}}_{\text{Bayes}} = \int \boldsymbol{\beta} \cdot \text{Dir}(\boldsymbol{\beta}|\boldsymbol{\eta} + \mathbf{n}) d\boldsymbol{\beta}$$

$$(\widehat{\boldsymbol{\beta}}_{\text{Bayes}})_i = \frac{\eta_i + n_i}{\sum_{j=1}^V (\eta_j + n_j)} = \frac{\eta_i + n_i}{\sum_{j=1}^V \eta_j + N}$$



Outline

- Sampling Method

- Monte Carlo

- Markov Chain Monte Carlo

- Metropolis-Hastings Algorithm

- Gibbs Sampling

- Sampling Method for LDA

- Bayesian Estimation

- Dirichlet-Multinomial Conjugate

- Gibbs Sampling for LDA

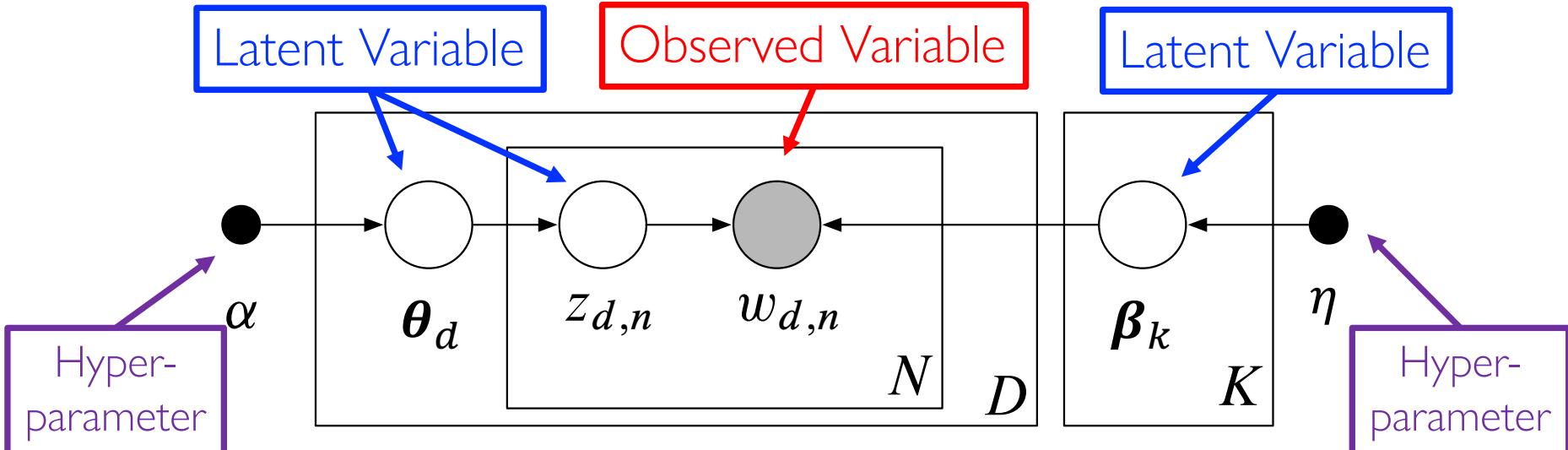


LDA as Latent Variable Model

- Two sets of random variables: z and x
 - z consists of **unobserved hidden variables**.
 - x consists of **observed variables**.
- Joint probability model parameterized by $\theta \in \Theta$:

$$p(x, z | \theta)$$

Griffiths et al. 2004



Gibbs Sampling for LDA

- We need to solve:

$$\text{In LDA: } p(\theta, z, \beta | w, \alpha, \eta) = \frac{p(\theta, z, \beta, w | \alpha, \eta)}{p(w | \alpha, \eta)}$$

- But the denominator is **intractable**!
- With Gibbs Sampling, we can **sample** from it!
- Shall we sample each vector in order?
 - $\theta_1, \theta_2, \dots, \theta_D, \beta_1, \beta_2, \dots, \beta_K, z_{11}, z_{12}, \dots, z_{dn}$.
- No! We only **sample** $\theta_{11}, \theta_{12}, \dots, \theta_{1K}$ in order:
 - Each dimension cannot be easily **disentangled** in Dirichlet variable!

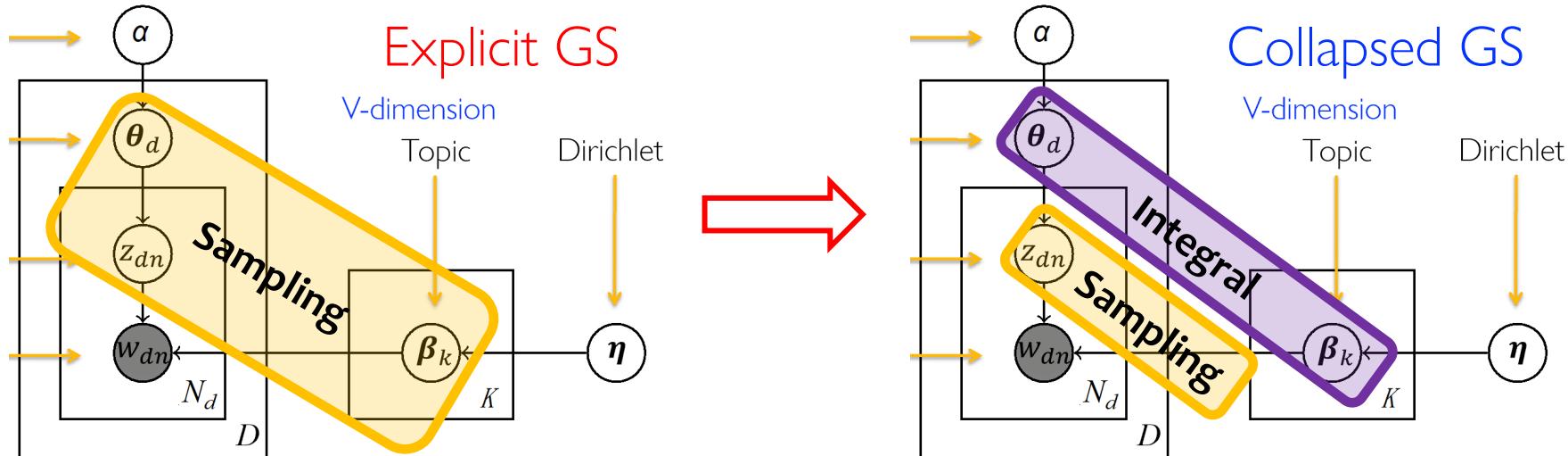
Gregor Heinrich. Parameter estimation for text analysis. Tech. Report.





Gibbs Sampling for LDA

- Sampling from a high-dimension distribution is very difficult.
- We need an **approximation** here:
 - Do not sample, but just compute mean of Posterior distribution.
- The question changes to: If we know all z_{dn} and w_{dn} (sampling!)
 - What are $p(\theta_d | z, w, \beta, \alpha, \eta)$, $p(\beta_k | z, w, \theta, \alpha, \eta)$ and their mean.

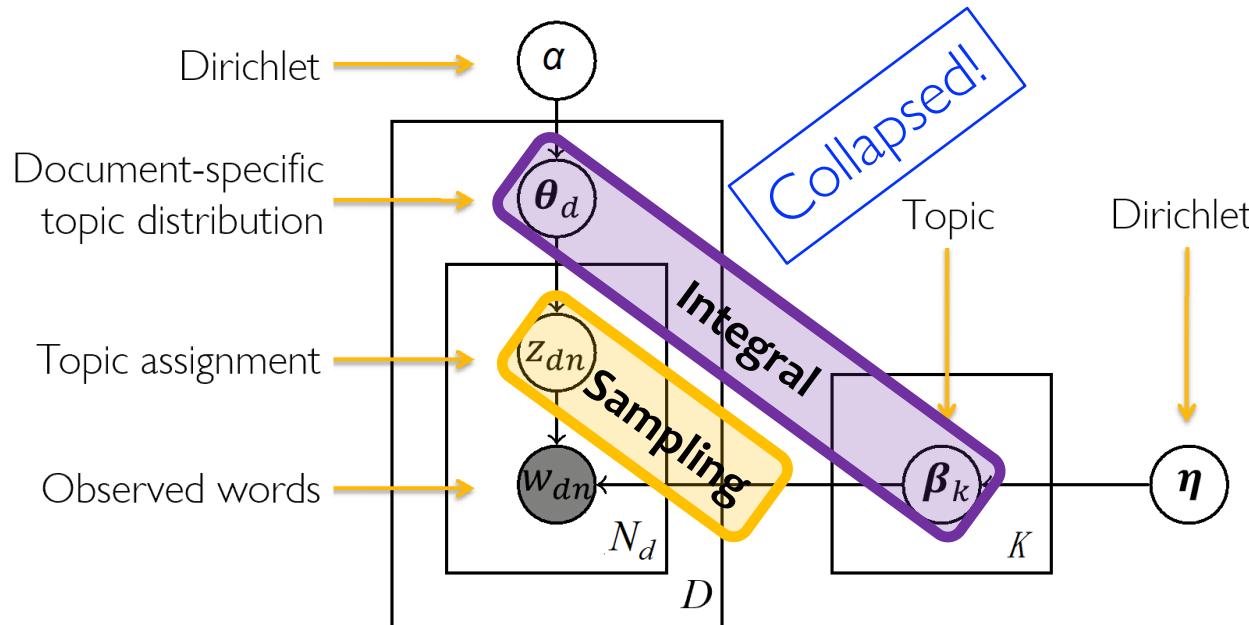




Collapsed Gibbs Sampling for LDA

- Iterate between sampling and integral:

1. Sample latent topic $z_{11}, z_{12}, \dots, z_{DN_d}$: $z_{ij} \sim p(z|z^{-ij}, w, \theta, \beta)$
2. Compute mean of $\theta_1, \theta_2, \dots, \theta_D$: $\widehat{\theta}_d = \mathbb{E}_{\theta_d \sim p(\theta_d|z, \alpha)}[\theta_d]$
3. Compute mean of $\beta_1, \beta_2, \dots, \beta_K$: $\widehat{\beta}_k = \mathbb{E}_{\beta_k \sim p(\beta_k|z, w, \eta)}[\beta_k]$



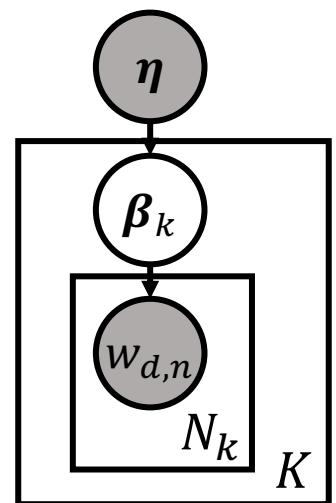
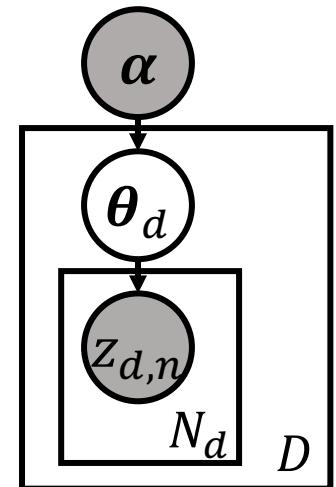
Terminology

- Component α :

- $N_d \leftarrow$ count of all words in document d
- $t_{dk} \leftarrow$ count of words in topic k in document d
- Count vector: $\mathbf{t}_d = (t_{d1}, \dots, t_{dK})$
- Conjugate in LDA: $p(\boldsymbol{\theta}_d = \boldsymbol{\theta} | \boldsymbol{\alpha}, \mathbf{z}) = \text{Dir}(\boldsymbol{\theta} | \boldsymbol{\alpha} + \mathbf{t}_d)$

- Component η :

- $N_k \leftarrow$ count of all words in topic k
- $n_{kv} \leftarrow$ count of words of term v in topic k
- Count vector: $\mathbf{n}_k = (n_{k1}, \dots, n_{kV})$
- Conjugate in LDA: $p(\boldsymbol{\beta}_k = \boldsymbol{\beta} | \boldsymbol{\eta}, \mathbf{w}, \mathbf{z}) = \text{Dir}(\boldsymbol{\beta} | \boldsymbol{\eta} + \mathbf{n}_k)$





Sampling from $p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}})$

- Gibbs sampling on latent topic from posterior: $p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}})$
- Problem: w_i is in the condition.
- We eliminate it with the Bayesian rule:

$$p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}) = p(z_i = k | w_i = v, \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}) = \frac{p(z_i = k, w_i = v | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i})}{p(w_i = v | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i})}$$

- The denominator is irrelevant to w_i :

$$p(w_i = v | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i})$$

- We have:

$$p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}) \propto p(z_i = k, w_i = v | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i})$$

- There is no w_i in the conditional distribution. Easy to sample!

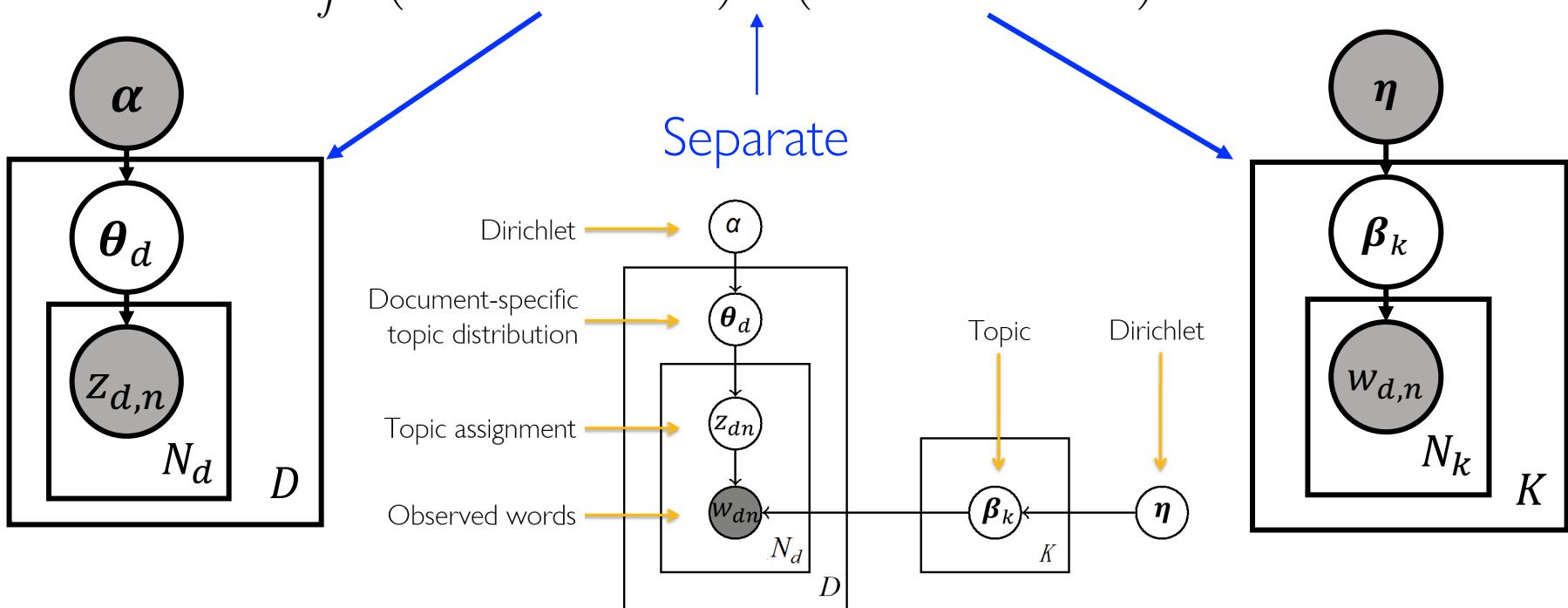


Sampling from $p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}})$

$$p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}) \propto p(z_i = k, w_i = v | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i})$$

$$= \int p(z_i = k, w_i = v, \vec{\theta}_d, \vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k$$

$$= \int p(z_i = k, \vec{\theta}_d | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) \cdot p(w_i = v, \vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k$$



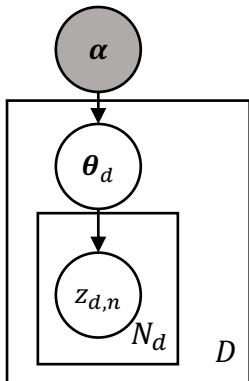
Gregor Heinrich. Parameter estimation for text analysis.





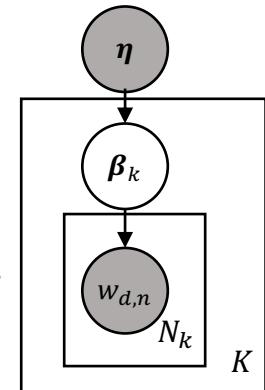
Sampling from $p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}})$

$$\begin{aligned}
 p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}) &\propto p(z_i = k, w_i = v | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) \\
 &= \int p(z_i = k, w_i = v, \vec{\theta}_d, \vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k \\
 &= \int p(z_i = k, \vec{\theta}_d | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) \cdot p(w_i = v, \vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k \\
 &= \int p(z_i = k | \vec{\theta}_d) p(\vec{\theta}_d | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) \cdot p(w_i = v | \vec{\beta}_k) p(\vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k \\
 &= \boxed{\int p(z_i = k | \vec{\theta}_d) \text{Dir}(\vec{\theta}_d | \vec{t}_{d,\neg i} + \vec{\alpha}) d\vec{\theta}_d \cdot \int p(w_i = v | \vec{\beta}_k) \text{Dir}(\vec{\beta}_k | \vec{n}_{k,\neg i} + \vec{\eta}) d\vec{\beta}_k}
 \end{aligned}$$



$$\rightarrow p(\boldsymbol{\theta}_d = \boldsymbol{\theta} | \boldsymbol{\alpha}, \mathbf{z}) = \text{Dir}(\boldsymbol{\theta} | \boldsymbol{\alpha} + \mathbf{t}_d)$$

$$p(\boldsymbol{\beta}_k = \boldsymbol{\beta} | \boldsymbol{\eta}, \mathbf{w}, \mathbf{z}) = \text{Dir}(\boldsymbol{\beta} | \boldsymbol{\eta} + \mathbf{n}_k) \leftarrow$$



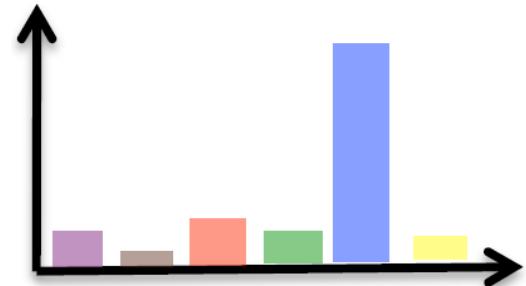
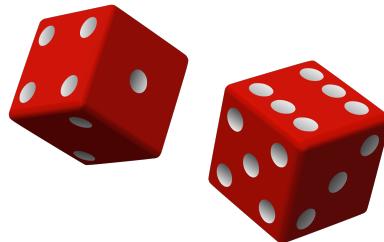
Gregor Heinrich. Parameter estimation for text analysis.





Sampling from $p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}})$

$$\begin{aligned} p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}) &\propto p(z_i = k, w_i = v | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) \\ &= \int p(z_i = k, w_i = v, \vec{\theta}_d, \vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k \\ &= \int p(z_i = k, \vec{\theta}_d | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) \cdot p(w_i = v, \vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k \\ &= \int p(z_i = k | \vec{\theta}_d) p(\vec{\theta}_d | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) \cdot p(w_i = v | \vec{\beta}_k) p(\vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k \\ &= \int p(z_i = k | \vec{\theta}_d) \text{Dir}(\vec{\theta}_d | \vec{t}_{d, \neg i} + \vec{\alpha}) d\vec{\theta}_d \cdot \int p(w_i = v | \vec{\beta}_k) \text{Dir}(\vec{\beta}_k | \vec{n}_{k, \neg i} + \vec{\eta}) d\vec{\beta}_k \\ &= \int \theta_{dk} \text{Dir}(\vec{\theta}_d | \vec{t}_{d, \neg i} + \vec{\alpha}) d\vec{\theta}_d \cdot \int \beta_{kv} \text{Dir}(\vec{\beta}_k | \vec{n}_{k, \neg i} + \vec{\eta}) d\vec{\beta}_k \\ &= \mathbb{E}_{\neg i}(\theta_{dk}) \cdot \mathbb{E}_{\neg i}(\beta_{kv}) \\ &= \hat{\theta}_{dk} \cdot \hat{\beta}_{kv} \end{aligned}$$



Gregor Heinrich. Parameter estimation for text analysis.





Sampling from $p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}})$

$$\begin{aligned}
 p(z_i = k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}) &\propto p(z_i = k, w_i = v | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) \\
 &= \int p(z_i = k, w_i = v, \vec{\theta}_d, \vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k \\
 &= \int p(z_i = k, \vec{\theta}_d | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) \cdot p(w_i = v, \vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k \\
 &= \int p(z_i = k | \vec{\theta}_d) p(\vec{\theta}_d | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) \cdot p(w_i = v | \vec{\beta}_k) p(\vec{\beta}_k | \vec{\mathbf{z}}_{\neg i}, \vec{\mathbf{w}}_{\neg i}) d\vec{\theta}_d d\vec{\beta}_k \\
 &= \int p(z_i = k | \vec{\theta}_d) \text{Dir}(\vec{\theta}_d | \vec{t}_{d, \neg i} + \vec{\alpha}) d\vec{\theta}_d \cdot \int p(w_i = v | \vec{\beta}_k) \text{Dir}(\vec{\beta}_k | \vec{n}_{k, \neg i} + \vec{\eta}) d\vec{\beta}_k \\
 &= \int \theta_{dk} \text{Dir}(\vec{\theta}_d | \vec{t}_{d, \neg i} + \vec{\alpha}) d\vec{\theta}_d \cdot \int \beta_{kv} \text{Dir}(\vec{\beta}_k | \vec{n}_{k, \neg i} + \vec{\eta}) d\vec{\beta}_k \\
 &= \mathbb{E}_{\neg i}(\theta_{dk}) \cdot \mathbb{E}_{\neg i}(\beta_{kv}) \\
 &= \hat{\theta}_{dk} \cdot \hat{\beta}_{kv}
 \end{aligned}$$

Bayes Estimator

$$\hat{\theta}_{dk, \neg i} = \frac{t_{dk, \neg i} + \alpha}{N_{d, \neg i} + \alpha K} \quad \hat{\beta}_{kv, \neg i} = \frac{n_{kv, \neg i} + \eta}{N_{k, \neg i} + \eta V}$$

Gregor Heinrich. Parameter estimation for text analysis.





Bayes Estimator for θ and β

- Document-topic distribution: $p(\theta_d = \theta | \alpha, z) = \text{Dir}(\theta | \alpha + t_d)$

$$\hat{\theta}_d = \int \theta \cdot \text{Dir}(\theta | \alpha + t_d) d\theta \quad \hat{\theta}_{dk} = \frac{t_{dk} + \alpha_k}{N_d + \sum_{k'=1}^K \alpha_{k'}}$$

- Term-topic distribution: $p(\beta_k = \beta | \eta, w, z) = \text{Dir}(\beta | \eta + n_k)$

$$\hat{\beta}_k = \int \beta \cdot \text{Dir}(\beta | \eta + n_k) d\beta \quad \hat{\beta}_{kv} = \frac{n_{kv} + \eta_v}{N_k + \sum_{v'=1}^V \eta_{v'}}$$

- Griffiths et al. 2004 further set $\alpha_k = \alpha$, $\eta_v = \eta$ as all constants.

$$\hat{\theta}_{dk} = \frac{t_{dk} + \alpha}{N_d + \alpha K} \quad \hat{\beta}_{kv} = \frac{n_{kv} + \eta}{N_k + \eta V}$$



LDA Inference with Gibbs Sampling

- Initialize all latent topics $\{z_{11}, z_{12}, \dots, z_{DN_D}\}$.

- Iterate:

- For each position ij , sample z_{ij} from posterior:

$$p(z_{ij} = k | \mathbf{z}_{\neg ij}, \mathbf{w}) \propto \frac{t_{dk, \neg ij} + \alpha}{N_{d, \neg ij} + \alpha K} \cdot \frac{n_{kv, \neg ij} + \eta}{N_{k, \neg ij} + \eta V}$$

- Compute Bayes Estimator of parameters:

$$\hat{\theta}_{dk} = \frac{t_{dk} + \alpha}{N_d + \alpha K} \quad \hat{\beta}_{kv} = \frac{n_{kv} + \eta}{N_k + \eta V}$$



Gregor Heinrich. Parameter estimation for text analysis.

Thank You Questions?

Mingsheng Long

mingsheng@tsinghua.edu.cn

<http://ise.thss.tsinghua.edu.cn/~mlong>

答疑：东主楼11区413室