

# Gibbs Sampling

**Subeen Cha**

Machine Learning Laboratory

Department of Statistics

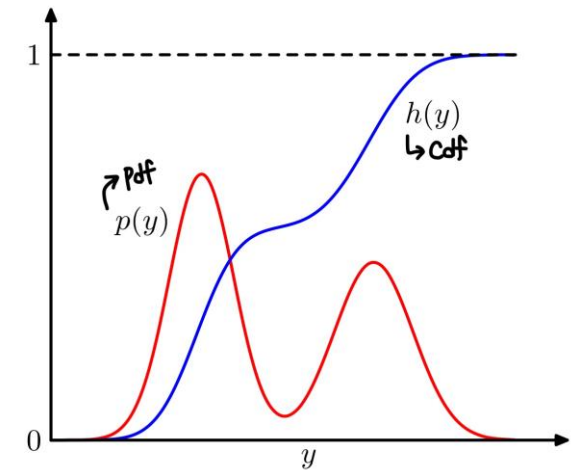
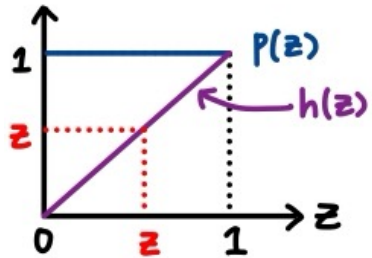
EWHA Womans University

- Sometimes, exact inference is intractable. → can use approximation
  - ex. Numerical sampling methods(Monte Carlo, etc)
- Sampling
  - General Idea: Obtain a set of samples  $\mathbf{z}^{(l)}$  ( $l = 1, \dots, L$ ) drawn independently from the distribution  $p(\mathbf{z})$ .
  - Then, the expectation can be approximated by a finite sum.
$$\mathbb{E}[\hat{f}] = \mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z} = \frac{1}{L} \sum_{l=1}^L f(\mathbf{z}^{(l)})$$
  - Some problems
    - The samples  $\{\mathbf{z}^{(l)}\}$  might not be independent.
    - The expectation may be dominated by regions of small probability.
      - $f(\mathbf{z})$ : small vs  $p(\mathbf{z})$ : large

## 1) Standard distributions

- Suppose that  $z$  is uniformly distributed over the interval  $(0, 1)$ , and that we transform the values of  $z$  using some function  $f(\cdot)$ .  $\Rightarrow y = f(z)$

$$z \sim U(0, 1), y = f(z) \rightarrow z = f^{-1}(y)$$



Jacobi transformation

$$h(y) \equiv F_Y(y) = P(Y \leq y) = \int_{-\infty}^y p(\hat{y}) \cdot d\hat{y} = P(f(z) \leq y) = P(z \leq f^{-1}(y)) = F_Z(f^{-1}(y)) = h(z) = z \quad (11.6)$$

$$p(y) \equiv f_Y(y) = \frac{d}{dy} F_Y(y) = \frac{d}{dy} F_Z(f^{-1}(y)) = f_Z(\underbrace{f^{-1}(y)}_z) \cdot \frac{d}{dy} \underbrace{f^{-1}(y)}_z = f_Z(z) \frac{dz}{dy} = p(z) \left| \frac{dz}{dy} \right| \quad (11.5)$$

$\rightarrow$  transform  $z$  using the inverse of the indefinite integral of the desired distribution  $p(y)$   
 $y = h^{-1}(z)$

## 1) Standard distributions

- Ex 1. exponential distribution

- pdf:  $p(y) = \lambda \exp(-\lambda y)$ ,  $0 \leq y < \infty$
- cdf:  $z = h(y) = \int_0^y p(\hat{y}) d\hat{y}$   
 $= 1 - \exp(-\lambda y)$   
 $\rightarrow \exp(-\lambda y) = 1 - z$   
 $\rightarrow -\lambda y = \ln(1 - z) \quad (\because 0 < z \leq 1)$   
 $\rightarrow y = -\lambda^{-1} \ln(1 - z) \sim \text{Exp}(\lambda)$

- Ex 2. Cauchy distribution

$$p(y) = \frac{1}{\pi} \cdot \frac{1}{1+y^2}$$

$$h(y) = \int_{-\infty}^y \frac{1}{\pi} \cdot \frac{1}{1+\hat{y}^2} d\hat{y} = \int_{-\frac{\pi}{2}}^{\tan^{-1}(y)} \frac{1}{\pi} \cdot \frac{1}{1+\tan^2 \theta} \cdot \sec^2 \theta d\theta$$

$\hat{y} = \tan \theta$   
 $d\hat{y} = \sec^2 \theta \cdot d\theta$

$$= \left[ \frac{1}{\pi} \cdot \theta \right]_{-\frac{\pi}{2}}^{\tan^{-1}(y)} = \frac{1}{\pi} \cdot \left\{ \tan^{-1}(y) + \frac{\pi}{2} \right\} = z$$

$$\rightarrow \tan^{-1}(y) = \pi z - \frac{\pi}{2}$$

$$\rightarrow y = \tan\left[\pi\left(z - \frac{1}{2}\right)\right]$$

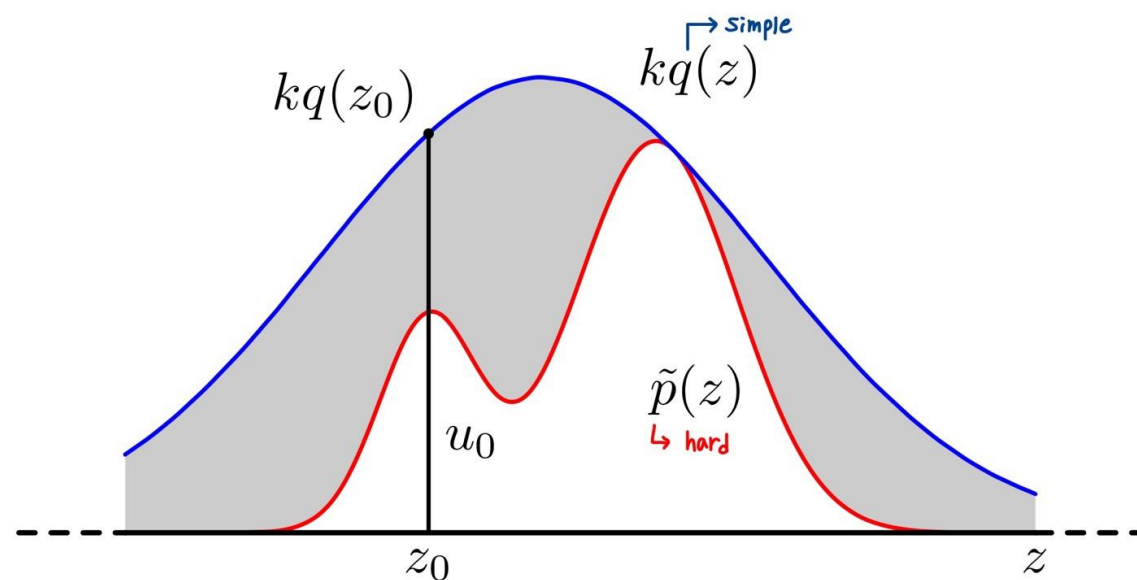
- multiple variables  $\rightarrow$  Use Jacobian

$$p(y_1, \dots, y_M) = p(z_1, \dots, z_M) \left| \frac{\partial(z_1, \dots, z_M)}{\partial(y_1, \dots, y_M)} \right|$$

- Generally, calculating and then inverting the indefinite integral of the required distribution is intractable.
  - More general strategy is required: rejection sampling and importance sampling

## 2) Rejection sampling

- Suppose that sampling directly from  $p(\mathbf{z})$  is difficult, but easily able to evaluate  $p(\mathbf{z})$  for any given value of  $\mathbf{z}$ :  $p(\mathbf{z}) = \frac{1}{Z_p} \tilde{p}(\mathbf{z})$ ,  $Z_p$ : unknown
- proposal distribution( $q(\mathbf{z})$ )
  - simpler distribution that we can readily draw samples
  - Choose a constant  $k$  s.t  $kq(\mathbf{z}) \geq \tilde{p}(\mathbf{z})$  for all values of  $\mathbf{z}$ .
    - $kq(\mathbf{z})$ : comparison function
- Each step of the rejection sampler involves generating two random numbers.
  - $z_0 \sim q(\mathbf{z})$
  - $u_0 \sim U[0, kq(z_0)]$
  - If  $u_0 > \tilde{p}(z_0)$ 
    - reject  $z_0$  and resample(grey shaded region)
  - Else
    - Accept  $z_0$
  - Samples are accepted with probability  $\tilde{p}(\mathbf{z})/kq(\mathbf{z})$ .
$$p(\text{accept}) = \int \{\tilde{p}(\mathbf{z})/kq(\mathbf{z})\}q(\mathbf{z})d\mathbf{z} = \frac{1}{k} \int \tilde{p}(\mathbf{z})d\mathbf{z}$$
$$\Rightarrow \text{choose } k \text{ as small as possible}$$



## 2- $\alpha$ ) Adaptive Rejection sampling

- Construction of an envelope function( $kq(z)$ ) is particularly easy if  $p(z)$  is log-concave.

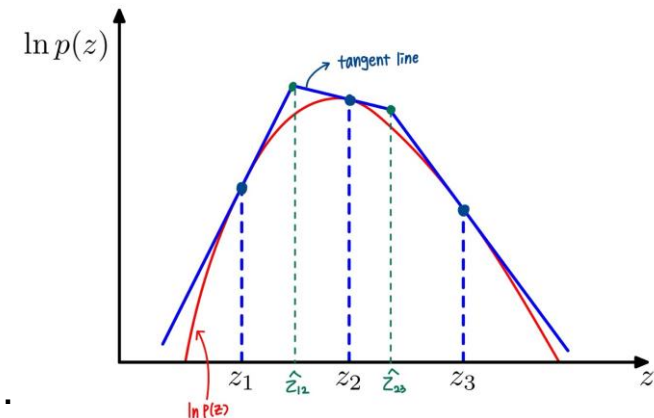
- Choose some initial set of grid points:  $z_1, z_2, z_3$
- Calculate  $\ln p(z_i)$  and gradient( $\lambda_i$ ).
- Construct an envelope function( $q(z)$ ).

- set of piecewise exponential distributions

$$q(z) = k_i \lambda_i \exp\{-\lambda_i(z - z_i)\}, \quad \widehat{z}_{i-1,i} < z \leq \widehat{z}_{i,i+1}, \quad k_i: \text{gradient of the tangent line at } z_i$$

- Once a sample has been drawn, apply the usual rejection criterion.

- If  $u_0 \leq \tilde{p}(z_0)$ 
  - Accept  $z_0$
- Else
  - Reject  $z_0$ .  $\rightarrow$  Incorporate into the set of grid points.
  - Compute a new tangent line.  $\rightarrow$  update envelope function  $\Rightarrow$  adaptive  
 $\Rightarrow$  Envelope function becomes a better approximation of the desired distribution  $p(z)$ .

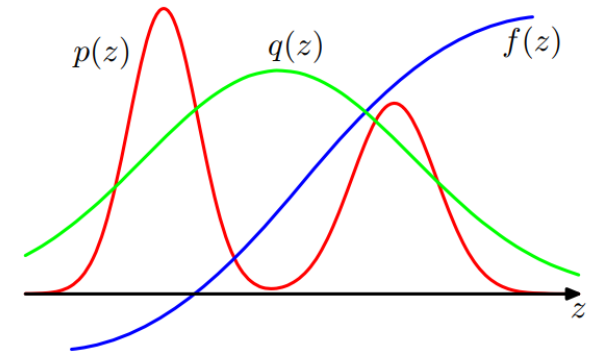


## 2) Rejection sampling\_limitation

- For high-dimensional problems, the acceptance rate decreases exponentially as the number of dimensions increases.
- Ex> multivariate Gaussian distribution
  - proposal distribution( $kq(z)$ ):  $N(0, \sigma_p^2 I)$
  - In  $D$ -dimensions the optimum value of  $k$  is given by  $k = (\sigma_q/\sigma_p)^D$ .
  - If the proposal distribution's variance is only slightly larger than the target distribution's, the acceptance ratio can become extremely low.
    - about 1/20,000 for  $D = 1,000$  dimensions
- Finding a good proposal distribution in high-dimensional problem is difficult.
  - Alternatives) importance sampling, etc

## 3) Importance sampling(= weighted sampling)

- Main goal of sampling
  - Generate the exact pdf. (X)
  - Calculating the expectation of pdf or calculating a certain probability (O)
- Uniform sampling
  - discretize  $\mathbf{z}$ -space into a uniform grid and evaluate partial sums
$$\mathbb{E}[f] \simeq \sum_{l=1}^L f(z^{(l)})p(z^{(l)})$$
  - But the number of terms in the summation grows exponentially with the dimensionality of  $\mathbf{z}$ .
  - In addition, most probability distributions have most masses inside relatively small regions in  $\mathbf{z}$  space in high-dimensional problems.
    - Only a very small proportion of the samples will make a significant contribution to the sum.



⇒ importance sampling



## 3) Importance sampling

- Again, use proposal distribution  $q(z)$  from which we can easily draw samples.
  - Sampling from  $q(z)$ :  $\{z^{(l)}\}$
  - Then, calculate the expectation.

$$\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z} = \int f(\mathbf{z})\frac{p(\mathbf{z})}{q(\mathbf{z})}q(\mathbf{z})d\mathbf{z} \simeq \frac{1}{L}\sum_{l=1}^L \frac{p(\mathbf{z}^{(l)})}{q(\mathbf{z}^{(l)})}f(\mathbf{z}^{(l)})$$

- $L$ : # of samples,  $\mathbf{z}^{(l)}$ : sample of  $Z$
- Importance weight:  $r_l = p(\mathbf{z}^{(l)})/q(\mathbf{z}^{(l)})$ 
  - kind of weighted average
  - used to correct the bias introduced by sampling from the wrong distribution( $q(z)$ )

## 3) Importance sampling

- Normalization

- case that the distribution  $p(\mathbf{z})$  can only be evaluated up to a normalization constant
  - i.e.,  $p(\mathbf{z}) = \tilde{p}(\mathbf{z})/Z_p$  where  $\tilde{p}(\mathbf{z})$  can be evaluated easily, whereas  $Z_p$  is unknown
- can use  $q(\mathbf{z}) = \tilde{q}(\mathbf{z})/Z_q$

$$\begin{aligned}\mathbb{E}[f] &= \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z} \\ &= \int f(\mathbf{z})\frac{p(\mathbf{z})}{q(\mathbf{z})}q(\mathbf{z})d\mathbf{z} = \int f(\mathbf{z})\frac{\tilde{p}(\mathbf{z})/Z_p}{\tilde{q}(\mathbf{z})/Z_q}q(\mathbf{z})d\mathbf{z} \\ &= \frac{Z_q}{Z_p} \int f(\mathbf{z})\frac{\tilde{p}(\mathbf{z})}{\tilde{q}(\mathbf{z})}q(\mathbf{z})d\mathbf{z} \\ &\simeq \frac{Z_q}{Z_p} \frac{1}{L} \sum_{l=1}^L \tilde{r}_l f(\mathbf{z}^{(l)}).\end{aligned}$$

- $\tilde{r}_l = \frac{\tilde{p}(\mathbf{z}^{(l)})}{\tilde{q}(\mathbf{z}^{(l)})}$
- $\frac{Z_p}{Z_q} = \frac{1}{Z_q} \int \tilde{p}(\mathbf{z})d\mathbf{z} = \int \frac{\tilde{p}(\mathbf{z})}{\tilde{q}(\mathbf{z})}\tilde{q}(\mathbf{z})d\mathbf{z} \simeq \frac{1}{L} \sum_{l=1}^L \tilde{r}_l$

$$\begin{aligned}\mathbb{E}[f] &\simeq \sum_{l=1}^L w_l f(\mathbf{z}^{(l)}) \\ \text{where } w_l &= \frac{\tilde{r}_l}{\sum_m \tilde{r}_m} = \frac{\tilde{p}(\mathbf{z}^{(l)})/q(\mathbf{z}^{(l)})}{\sum_m \tilde{p}(\mathbf{z}^{(m)})/q(\mathbf{z}^{(m)})}\end{aligned}$$

## 3) Importance sampling\_likelihood weighted sampling

- Significant difference from  $p(z)$  and  $q(z)$  can reduce the efficiency of sampling.
  - improvement  $\rightarrow$  likelihood weighted sampling
- Likelihood weighted sampling
  - based on ancestral sampling of the variables
    - make one pass through the set of variables in the order  $z_1, \dots, z_M$  sampling from the conditional distributions  $p(\mathbf{z}_i | \text{pa}_i)$ .
  - Check each variable in turn.
    - If that variable is in the evidence set(= observed variables), then it is just set to its instantiated value.
    - Else, sampling from the conditional distribution  $p(\mathbf{z}_i | \text{pa}_i)$
  - Weights

$$r(\mathbf{z}) = \prod_{\mathbf{z}_i \notin \mathbf{e}} \frac{p(\mathbf{z}_i | \text{pa}_i)}{p(\mathbf{z}_i | \text{pa}_i)} \prod_{\mathbf{z}_i \in \mathbf{e}} \frac{p(\mathbf{z}_i | \text{pa}_i)}{1} = \prod_{\mathbf{z}_i \in \mathbf{e}} p(\mathbf{z}_i | \text{pa}_i)$$

## 3- $\alpha$ ) Sampling Importance Resampling

Improves weight balance problem in importance sampling

- Process

- 1. Sampling: draw  $L$  samples  $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(L)}$  from  $q(\mathbf{z})$ .
- 2. Weight construction:  $w_l = \frac{\tilde{r}_l}{\sum_m \tilde{r}_m} = \frac{\tilde{p}(\mathbf{z}^{(l)})/q(\mathbf{z}^{(l)})}{\sum_m \tilde{p}(\mathbf{z}^{(m)})/q(\mathbf{z}^{(m)})}$
- 3. Resampling: resample #  $L$  samples from the sample sets  $(\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(L)})$ .
  - Sampling probabilities are given by the weights  $(w_1, \dots, w_L)$ .
- The approximation improves as the sampling distribution  $q(\mathbf{z})$  gets closer to the desired distribution  $p(\mathbf{z})$ .
  - If  $L \rightarrow \infty$ , the distribution becomes correct.
    - When  $q(\mathbf{z}) = p(\mathbf{z})$ , the samples sets  $(\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(L)})$  have the desired distribution  $(p(\mathbf{z}))$ , and the weights become uniformly  $w_n = 1/L$ .
    - So, the resampled values also have the desired distribution.

## Sampling and the EM Algorithm

- M-step

- optimize  $Q(\log\text{-likelihood estimates})$  w.r.t  $\theta$  (model parameters).
- Update the model parameter  $\theta$  given the value  $z$ .

$$Q(\theta, \theta^{old}) = \int p(\mathbf{Z}|\mathbf{X}, \theta^{old}) \ln p(\mathbf{Z}, \mathbf{X}|\theta) d\mathbf{Z}$$

- E-step

- Update  $z$  given the model parameter value  $\theta$ .
- can use sampling methods to approximate the integral by a finite sum over samples  $\{\mathbf{Z}^{(l)}\}$

$$Q(\theta, \theta^{old}) \simeq \frac{1}{L} \sum_{l=1}^L \ln p(\mathbf{Z}^{(l)}, \mathbf{X}|\theta)$$

## 2-1. Markov Chain

- A Markov chain(or Markov process) is a stochastic model describing a sequence of possible events in which the probability of each event depends only on the state attained in the previous event. © Wikipedia
- A series of random variables  $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(M)}$  s.t the following conditional independence property holds
  - Conditional independence:  $p(\mathbf{z}^{(m+1)} | \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}) = p(\mathbf{z}^{(m+1)} | \mathbf{z}^{(m)}), \quad m \in \{1, \dots, M - 1\}$



## 2-1. Markov Chain

- Terminologies & Properties

- transition probabilities

- We can specify the Markov chain by giving the probability distribution for the initial variable  $p(\mathbf{z}^{(0)})$  together with the conditional probabilities for subsequent variables.

$$T_m(z^{(m)}, z^{(m+1)}) \equiv p(z^{(m+1)} | z^{(m)})$$

- A Markov chain is called **homogeneous** if the transition probabilities are the same for all  $m$ .

- invariant(= stationary)

- A distribution is said to be invariant(or stationary) if each step in the chain leaves that distribution invariant.
    - For a homogeneous Markov chain with transition probabilities  $T(z', z)$ , the distribution  $p^*(z)$  is invariant if  $p^*(z) = \sum_{z'} T(z', z)p^*(z')$ .

## 2-1. Markov Chain

- Terminologies & Properties

- detailed balance condition

$$p^*(\mathbf{z})T(\mathbf{z}, \mathbf{z}') = p^*(\mathbf{z}')T(\mathbf{z}', \mathbf{z})$$

- Transition probability that satisfies detailed balance w.r.t a particular distribution will leave that distribution invariant.

$$\sum_{\mathbf{z}'} p^*(\mathbf{z}')T(\mathbf{z}', \mathbf{z}) = \sum_{\mathbf{z}'} p^*(\mathbf{z})T(\mathbf{z}, \mathbf{z}') = p^*(\mathbf{z}) \sum_{\mathbf{z}'} p(\mathbf{z}'|\mathbf{z}) = p^*(\mathbf{z}).$$

- A Markov chain that respects detailed balance is said to be **reversible**.

- Ergodicity

- For  $m \rightarrow \infty$ , the distribution  $p(\mathbf{z}^{(m)})$  converges to the required invariant distribution  $p^*(\mathbf{z})$ , irrespective of the choice of initial distribution  $p(\mathbf{z}^{(0)})$ .
- $p^*(\mathbf{z})$ : equilibrium distribution



## 2-1. Markov Chain

- Markov Chain for Sampling
  - In previous, we do not use the past records.
    - Every sampling is independent.
  - MCMC generates continuous samples.
    - $\{z^{(1)}, z^{(2)}, \dots\}$ : forms a Markov Chain
    - Also, maintain a record of the current state  $z^{(\tau)}$  and the proposal distribution  $q(z|z^{(\tau)})$ .
    - At each cycle of the algorithm, we generate a candidate sample  $z^*$  from the proposal distribution and then accept the sample according to an appropriate criterion.

## 2-2. The Metropolis-Hastings Algorithm

- Metropolis Algorithm

General algorithm of MCMC

- Assumption) The proposal distribution is symmetric:  $q(\mathbf{z}_A|\mathbf{z}_B) = q(\mathbf{z}_B|\mathbf{z}_A)$
- Process
  - Current value:  $z^{(\tau)}$
  - Propose a candidate  $z^* \sim q(z^*|z^{(\tau)})$ .  $q(\cdot)$  : proposal distribution
  - With an acceptance probability  $A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) = \min\left(1, \frac{\tilde{p}(\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(\tau)})}\right)$ ,
    - Choose a random number  $u$  with uniform distribution over the unit interval  $(0, 1)$ .
    - If  $A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) > u$ : accept,  $\mathbf{z}^{(\tau+1)} \rightarrow \mathbf{z}^*$
    - Else: reject,  $\mathbf{z}^{(\tau+1)} \rightarrow \mathbf{z}^{(\tau)}$ 
      - The candidate point  $\mathbf{z}^*$  is discarded. + The previous sample is included instead in the final list of samples.
      - Another candidate sample is drawn from the distribution.

## 2-2. The Metropolis-Hastings Algorithm

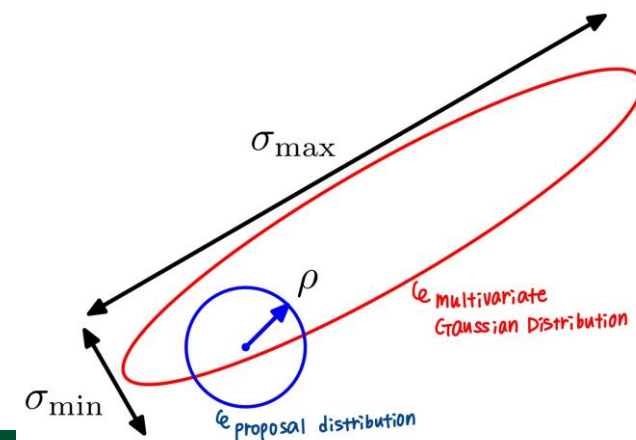
- Metropolis-Hastings Algorithm

- Remove the symmetric assumption. → more general case
- Now, acceptance probability becomes  $A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) = \min \left( 1, \frac{\tilde{p}(\mathbf{z}^*) q_k(\mathbf{z}^{(\tau)} | \mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(\tau)}) q_k(\mathbf{z}^* | \mathbf{z}^{(\tau)})} \right)$ .
- $p(\mathbf{z})$  is an invariant distribution of the Markov chain defined by the Metropolis-Hastings algorithm.
  - Pf> Show that detailed balance condition holds.
    - Detailed balance condition:  $p^*(\mathbf{z})T(\mathbf{z}, \mathbf{z}') = p^*(\mathbf{z}')T(\mathbf{z}', \mathbf{z})$

$$\begin{aligned} p(\mathbf{z}) q_k(\mathbf{z}' | \mathbf{z}) A_k(\mathbf{z}', \mathbf{z}) &= \min (p(\mathbf{z}) q_k(\mathbf{z}' | \mathbf{z}), p(\mathbf{z}') q_k(\mathbf{z} | \mathbf{z}')) \\ &= \min (p(\mathbf{z}') q_k(\mathbf{z} | \mathbf{z}'), p(\mathbf{z}) q_k(\mathbf{z}' | \mathbf{z})) \\ &= p(\mathbf{z}') q_k(\mathbf{z} | \mathbf{z}') A_k(\mathbf{z}, \mathbf{z}') \end{aligned}$$

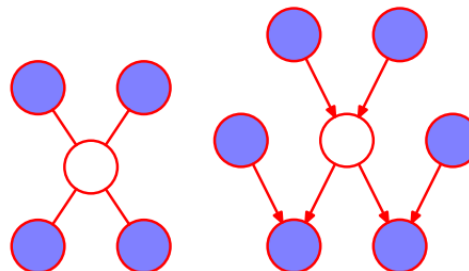
## 2-2. The Metropolis-Hastings Algorithm

- Metropolis-Hastings Algorithm
  - The specific choice of proposal distribution can have a marked effect on the performance of the algorithm.
    - continuous state spaces: Gaussian centered on the current state
  - There is an important trade-off in determining the variance parameter.
    - Small variance: high transition approval rate, but slow progress through the state space
    - Large variance: high rejection rate
  - The size of the proposed distribution( $\rho$ ) should be of the same order as the smallest length scale( $\sigma_{\min}$ ).
    - explores the distribution along the more extended direction by means of a random walk
    - can have very slow convergence if the distributions vary are very different in different directions



## Intro

- Gibbs Sampling: special case of the Metropolis-Hastings algorithm
- Consider a Metropolis-Hastings sampling step involving the variable  $z_k$ .
  - $\mathbf{z}^{(\tau)} = (z_k^{(\tau)}, \mathbf{z}_{\setminus k}^{(\tau)}) \rightarrow \mathbf{z}^* = (z_k^*, \mathbf{z}_{\setminus k}^*)$
  - the remaining variables  $\mathbf{z}_{\setminus k}$  remain fixed:  $\mathbf{z}_{\setminus k}^* = \mathbf{z}_{\setminus k}$
  - transition probability:  $q_k(\mathbf{z}^*|\mathbf{z}) = p(z_k^*|\mathbf{z}_{\setminus k})$
  - Then, the acceptance probability becomes 1.  $\rightarrow$  always accepted
    - $p(\mathbf{z}) = p(z_k|\mathbf{z}_{\setminus k})p(\mathbf{z}_{\setminus k}) \rightarrow A(\mathbf{z}^*, \mathbf{z}) = \frac{p(\mathbf{z}^*)q_k(\mathbf{z}|\mathbf{z}^*)}{p(\mathbf{z})q_k(\mathbf{z}^*|\mathbf{z})} = \frac{p(z_k^*|\mathbf{z}_{\setminus k}^*)p(\mathbf{z}_{\setminus k}^*)p(z_k|\mathbf{z}_{\setminus k}^*)}{p(z_k|\mathbf{z}_{\setminus k})p(\mathbf{z}_{\setminus k})p(z_k^*|\mathbf{z}_{\setminus k})} = 1$
  - Applicability
    - depends on the ease with which samples can be drawn from the conditional distributions  $p(z_k|\mathbf{z}_{\setminus k})$
    - Can simplify using the Markov blanket



## Concept of Gibbs Sampling

- Each step involves replacing the value of one of the variables by a value drawn from the distribution of that variable conditioned on the values of the remaining variables.
- Repeated either by cycling through the variables in some particular order or by choosing the variable to be updated at each step at random from some distribution
- Process
  1. Initialize  $\{z_i : i = 1, \dots, M\}$
  2. For  $\tau = 1, \dots, T$ :
    - Sample  $z_1^{(\tau+1)} \sim p(z_1 | z_2^{(\tau)}, z_3^{(\tau)}, \dots, z_M^{(\tau)})$ .
    - Sample  $z_2^{(\tau+1)} \sim p(z_2 | z_1^{(\tau+1)}, z_3^{(\tau)}, \dots, z_M^{(\tau)})$ .
    - $\vdots$
    - Sample  $z_j^{(\tau+1)} \sim p(z_j | z_1^{(\tau+1)}, \dots, z_{j-1}^{(\tau+1)}, z_{j+1}^{(\tau)}, \dots, z_M^{(\tau)})$ .
    - $\vdots$
    - Sample  $z_M^{(\tau+1)} \sim p(z_M | z_1^{(\tau+1)}, z_2^{(\tau+1)}, \dots, z_{M-1}^{(\tau+1)})$ .

## Intro

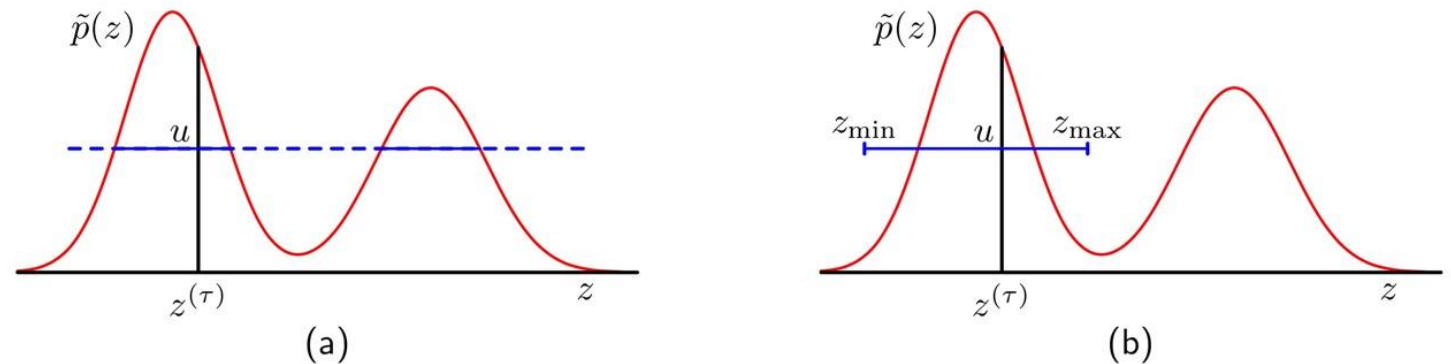
- Metropolis algorithm is sensitive to the step size.
  - too small: slow decorrelation due to random walk behavior
  - too large: inefficiency due to a high rejection rate

⇒ Slice Sampling provides an adaptive step size that is automatically adjusted to match the characteristics of the distribution.
- Univariate case
  - Introduce an auxiliary variable  $u$  and draw samples from the joint  $(z, u)$  space.
  - Goal: sample uniformly from the area under the distribution
    - joint distribution: 
$$\hat{p}(z, u) = \begin{cases} 1/Z_p & \text{if } 0 \leq u \leq \tilde{p}(z) \\ 0 & \text{otherwise} \end{cases}$$
    - marginal distribution over  $z$ : 
$$p(z) = \int \tilde{p}(z, u) du = \tilde{p}(z)/Z_p$$

## Methology

- Procedure

- Alternately sample  $z$  and  $u$ .
  - Given the current value  $z^{(\tau)}$ , select  $u$  uniformly in the range  $0 \leq u \leq \tilde{p}(z)$ 
    - creates a "slice" through the distribution
  - Then, fix  $u$  and sample  $z$  from the 'slice' through the distribution defined by  $\{z: \tilde{p}(z) > u\}$
- Adjust the region based on the characteristic length scales of the distribution
- Repeat until a valid  $z$  is found.



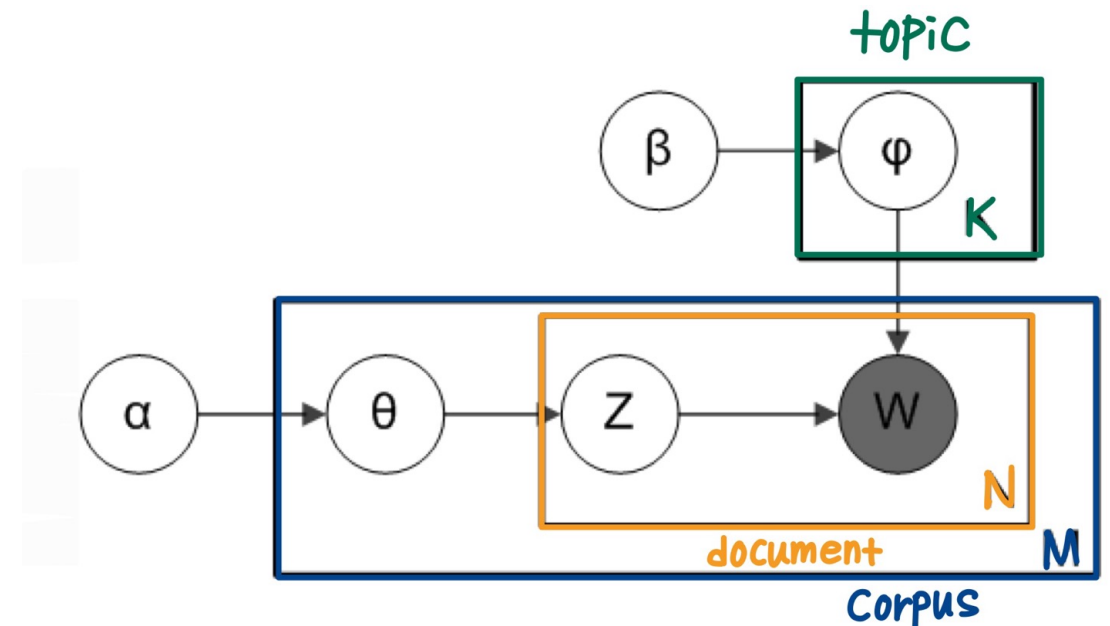
**Figure 11.13** Illustration of slice sampling. (a) For a given value  $z^{(\tau)}$ , a value of  $u$  is chosen uniformly in the region  $0 \leq u \leq \tilde{p}(z^{(\tau)})$ , which then defines a 'slice' through the distribution, shown by the solid horizontal lines. (b) Because it is infeasible to sample directly from a slice, a new sample of  $z$  is drawn from a region  $z_{\min} \leq z \leq z_{\max}$ , which contains the previous value  $z^{(\tau)}$ .



# LATENT DIRICHLET ALLOCATION

## Review

- Graphical model
  - $M$ : # of documents - corpus
  - $N$ : # of words in a given document (document  $i$  has  $N_i$  words)
  - $\alpha$ : parameter of the Dirichlet prior on the per-document topic distributions( $\theta$ )
  - $\beta$ : parameter of the Dirichlet prior on the per-topic word distribution( $\varphi$ )
  - $\theta_i$ : topic distribution for document  $i$
  - $\varphi_k$ : word distribution for topic  $k$
  - $z_{ij}$ : topic for the  $j$ -th word in document  $i$
  - $w_{ij}$ : specific word

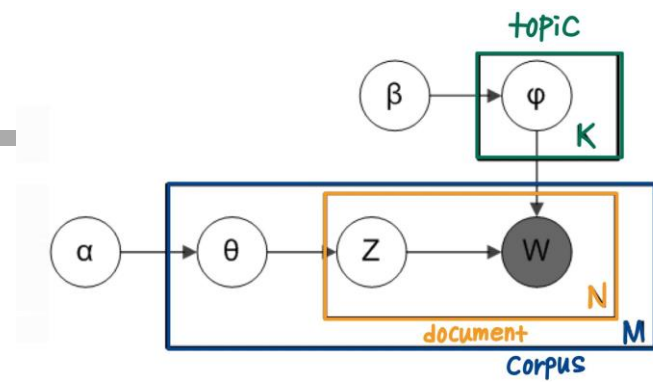


# LDA – Gibbs Sampling

## Review

- Generative process

- $\theta_i \sim \text{Dir}(\alpha), i \in \{1, \dots, M\}$
  - $\varphi_k \sim \text{Dir}(\beta), k \in \{1, \dots, K\}$
  - $z_{ij} \sim \text{Multinomial}(\theta_i), i \in \{1, \dots, M\}, j \in \{1, \dots, N\}$
  - $w_{ij} \sim \text{Multinomial}(\varphi_{z_{ij}}), i \in \{1, \dots, M\}, j \in \{1, \dots, N\}$
- A word  $w$  is generated from the distribution of  $\varphi_z$  word-topic distribution.
- $z$  topic is generated from the distribution of  $\theta$  document-topic distribution.
- $\theta$  document topic distribution is generated from the distribution of  $\alpha$ .
- $\varphi$  word-topic distribution is generated from the distribution of  $\beta$ .
- We want to find the most likely allocation of  $\mathbf{Z}$ .
  - If we have  $\mathbf{Z}$  distribution, we can find the most likely  $\theta$  and  $\varphi$ .

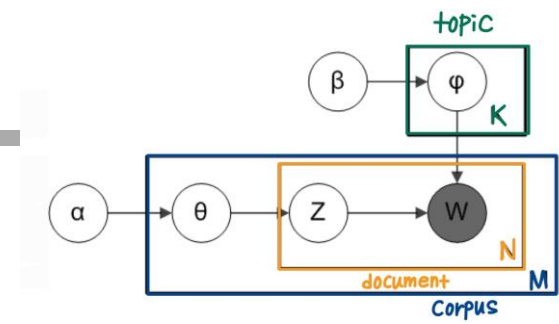


# LDA – Gibbs Sampling

## Collapsed Gibbs Sampling

Find the most likely assignment on  $\mathbf{Z}$ .

- Start with the factorization.  $P(\mathbf{W}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\varphi}; \alpha, \beta) = \prod_{i=1}^K P(\varphi_i; \beta) \prod_{j=1}^M P(\theta_j; \alpha) \prod_{t=1}^N P(Z_{j,t} | \theta_j) P(W_{j,t} | \varphi_{Z_{j,t}}),$
- Collapse  $\theta$  and  $\varphi$ .
  - Leave only  $\mathbf{W}$ (data point),  $\mathbf{Z}$ (sampling target),  $\alpha, \beta$ (priors)
  - Marginalize out!



$$P(\mathbf{Z}, \mathbf{W}; \alpha, \beta) \stackrel{\text{marginalize}}{=} \int_{\boldsymbol{\theta}} \int_{\boldsymbol{\varphi}} P(\mathbf{W}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\varphi}; \alpha, \beta) d\boldsymbol{\varphi} d\boldsymbol{\theta}$$

$$\stackrel{\substack{\varphi \text{ and } \theta \\ \text{are independent}}}{=} \int_{\boldsymbol{\varphi}} \prod_{i=1}^K P(\varphi_i; \beta) \prod_{j=1}^M \prod_{t=1}^N P(W_{j,t} | \varphi_{Z_{j,t}}) d\boldsymbol{\varphi} \int_{\boldsymbol{\theta}} \prod_{j=1}^M P(\theta_j; \alpha) \prod_{t=1}^N P(Z_{j,t} | \theta_j) d\boldsymbol{\theta}.$$

# LDA – Gibbs Sampling

## Collapsed Gibbs Sampling

Find the most likely assignment on  $Z$ .

- Collapse  $\theta$  and  $\varphi$ .

- $\theta$  part

$$\int_{\varphi} \prod_{i=1}^K P(\varphi_i; \beta) \prod_{j=1}^M \prod_{t=1}^N P(W_{j,t} | \varphi_{Z_{j,t}}) d\varphi \int_{\theta} \prod_{j=1}^M P(\theta_j; \alpha) \prod_{t=1}^N P(Z_{j,t} | \theta_j) d\theta.$$

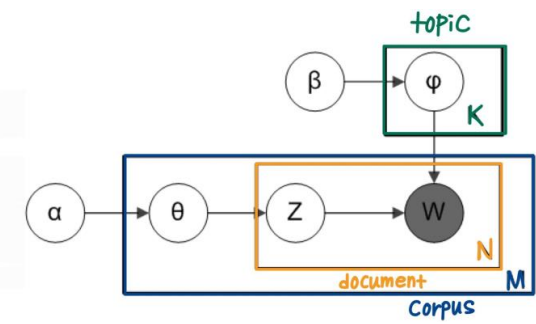
$$\int_{\theta} \prod_{j=1}^M P(\theta_j; \alpha) \prod_{t=1}^N P(Z_{j,t} | \theta_j) d\theta = \prod_{j=1}^M \int_{\theta_j} \left\{ P(\theta_j; \alpha) \prod_{t=1}^N P(Z_{j,t} | \theta_j) \right\} d\theta_j.$$

$$\int_{\theta_j} \frac{P(\theta_j; \alpha) \prod_{t=1}^N P(Z_{j,t} | \theta_j) d\theta_j}{\theta_j \sim \text{Dir}(\alpha)} = \int_{\theta_j} \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{j,i}^{\alpha_i-1} \prod_{t=1}^N P(Z_{j,t} | \theta_j) d\theta_j.$$

- $n_{j,r}^i$ : the number of word tokens in the  $j$ th document with the same word symbol (the  $r$ th word in the (total) vocabulary) assigned to the  $i$ th topic

$$P(Z_{j,t} = i | \theta_j) = \theta_{j,i} \rightarrow \prod_{t=1}^N P(Z_{j,t} | \theta_j) = \prod_{i=1}^K \theta_{j,i}^{n_{j,(\cdot)}^i}.$$

$$\int_{\theta_j} \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{j,i}^{\alpha_i-1} \prod_{i=1}^K \theta_{j,i}^{n_{j,(\cdot)}^i} d\theta_j = \int_{\theta_j} \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{j,i}^{n_{j,(\cdot)}^i + \alpha_i - 1} d\theta_j.$$



# LDA – Gibbs Sampling

## Collapsed Gibbs Sampling

Find the most likely assignment on  $Z$ .

- Collapse  $\theta$  and  $\varphi$ .

- $\theta$  part

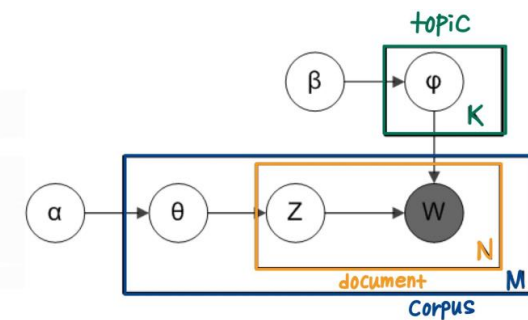
$$\int_{\theta_j} \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{j,i}^{\alpha_i-1} \prod_{i=1}^K \theta_{j,i}^{n_{j,i}^i} d\theta_j = \int_{\theta_j} \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{j,i}^{n_{j,i}^i + \alpha_i - 1} d\theta_j.$$

- Use pdf property:  $\int_{\theta_j} \frac{\Gamma(\sum_{i=1}^K n_{j,i}^i + \alpha_i)}{\prod_{i=1}^K \Gamma(n_{j,i}^i + \alpha_i)} \prod_{i=1}^K \theta_{j,i}^{n_{j,i}^i + \alpha_i - 1} d\theta_j \stackrel{\text{pdf}}{=} 1.$

$$p(P = \{p_i\} | \alpha_i) = \frac{\prod_i \Gamma(\alpha_i)}{\Gamma(\sum_i \alpha_i)} \prod_i p_i^{\alpha_i - 1}$$

- Make Dirichlet Distribution form:

$$\begin{aligned} \int_{\theta_j} \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{t=1}^N P(Z_{j,t} | \theta_j) d\theta_j &= \int_{\theta_j} \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{j,i}^{n_{j,i}^i + \alpha_i - 1} d\theta_j \\ &= \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \frac{\prod_{i=1}^K \Gamma(n_{j,i}^i + \alpha_i)}{\Gamma(\sum_{i=1}^K n_{j,i}^i + \alpha_i)} \int_{\theta_j} \frac{\Gamma(\sum_{i=1}^K n_{j,i}^i + \alpha_i)}{\prod_{i=1}^K \Gamma(n_{j,i}^i + \alpha_i)} \prod_{i=1}^K \theta_{j,i}^{n_{j,i}^i + \alpha_i - 1} d\theta_j \\ &= \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \frac{\prod_{i=1}^K \Gamma(n_{j,i}^i + \alpha_i)}{\Gamma(\sum_{i=1}^K n_{j,i}^i + \alpha_i)} \cdot \underbrace{\int_{\theta_j} \frac{\Gamma(\sum_{i=1}^K n_{j,i}^i + \alpha_i)}{\prod_{i=1}^K \Gamma(n_{j,i}^i + \alpha_i)} \prod_{i=1}^K \theta_{j,i}^{n_{j,i}^i + \alpha_i - 1} d\theta_j}_{\sim \text{Dir}(n_{j,\cdot}^i + \alpha_i)} \end{aligned}$$



# LDA – Gibbs Sampling

## Collapsed Gibbs Sampling

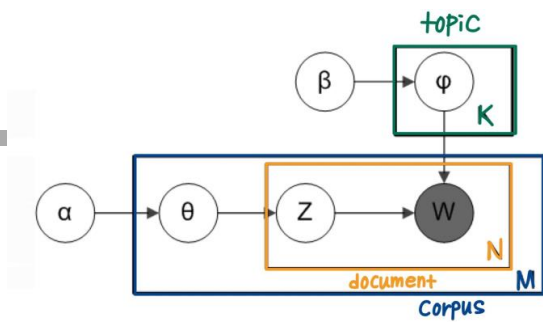
Find the most likely assignment on  $Z$ .

- Collapse  $\theta$  and  $\varphi$ .  $\int_{\varphi} \prod_{i=1}^K P(\varphi_i; \beta) \prod_{j=1}^M \prod_{t=1}^N P(W_{j,t} | \varphi_{Z_{j,t}}) d\varphi \int_{\theta} \prod_{j=1}^M P(\theta_j; \alpha) \prod_{t=1}^N P(Z_{j,t} | \theta_j) d\theta$ .

- $\varphi$  part

$$\begin{aligned}
 & \int_{\varphi} \prod_{i=1}^K P(\varphi_i; \beta) \prod_{j=1}^M \prod_{t=1}^N P(W_{j,t} | \varphi_{Z_{j,t}}) d\varphi \\
 &= \prod_{i=1}^K \int_{\varphi_i} \underline{P(\varphi_i; \beta)} \prod_{j=1}^M \prod_{t=1}^N P(W_{j,t} | \varphi_{Z_{j,t}}) d\varphi_i \\
 & \stackrel{\varphi \sim \text{Dir}(\beta)}{=} \prod_{i=1}^K \int_{\varphi_i} \frac{\Gamma(\sum_{r=1}^V \beta_r)}{\prod_{r=1}^V \Gamma(\beta_r)} \prod_{r=1}^V \varphi_{i,r}^{\beta_r-1} \prod_{r=1}^V \varphi_{i,r}^{n_{(\cdot),r}^i} d\varphi_i \\
 &= \prod_{i=1}^K \int_{\varphi_i} \frac{\Gamma(\sum_{r=1}^V \beta_r)}{\prod_{r=1}^V \Gamma(\beta_r)} \prod_{r=1}^V \varphi_{i,r}^{n_{(\cdot),r}^i + \beta_r - 1} d\varphi_i \\
 & \stackrel{\text{pdf property}}{=} \prod_{i=1}^K \frac{\Gamma(\sum_{r=1}^V \beta_r)}{\prod_{r=1}^V \Gamma(\beta_r)} \frac{\prod_{r=1}^V \Gamma(n_{(\cdot),r}^i + \beta_r)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^i + \beta_r)}.
 \end{aligned}$$

$$P(\mathbf{Z}, \mathbf{W}; \alpha, \beta) = \prod_{j=1}^M \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \frac{\prod_{i=1}^K \Gamma(n_{j,(\cdot)}^i + \alpha_i)}{\Gamma(\sum_{i=1}^K n_{j,(\cdot)}^i + \alpha_i)} \times \prod_{i=1}^K \frac{\Gamma(\sum_{r=1}^V \beta_r)}{\prod_{r=1}^V \Gamma(\beta_r)} \frac{\prod_{r=1}^V \Gamma(n_{(\cdot),r}^i + \beta_r)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^i + \beta_r)}.$$



# LDA – Gibbs Sampling

# Collapsed Gibbs Sampling

Find the most likely assignment on  $Z$ .

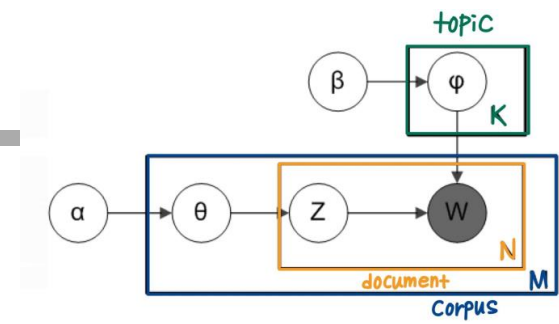
- Gibbs sampling formula

$$P(\mathbf{Z}, \mathbf{W}; \alpha, \beta) = \prod_{j=1}^M \frac{\Gamma\left(\sum_{i=1}^K \alpha_i\right)}{\prod_{i=1}^K \Gamma(\alpha_i)} \frac{\prod_{i=1}^K \Gamma(n_{j,(\cdot)}^i + \alpha_i)}{\Gamma\left(\sum_{i=1}^K n_{j,(\cdot)}^i + \alpha_i\right)} \times \prod_{i=1}^K \frac{\Gamma\left(\sum_{r=1}^V \beta_r\right)}{\prod_{r=1}^V \Gamma(\beta_r)} \frac{\prod_{r=1}^V \Gamma(n_{(\cdot),r}^i + \beta_r)}{\Gamma\left(\sum_{r=1}^V n_{(\cdot),r}^i + \beta_r\right)}.$$

- $W$ ,  $\alpha$ , and  $\beta$  are assumed or data points, and  $\mathbf{Z}$  is the target of sampling.
  - In Gibbs Sampling, we sample  $\mathbf{Z}$  one by one.
  - The key point is to derive the following conditional probability.

$$P(Z_{(m,n)} \mid \mathbf{Z}_{-(m,n)}, \mathbf{W}; \alpha, \beta) = \frac{P(Z_{(m,n)}, \mathbf{Z}_{-(m,n)}, \mathbf{W}; \alpha, \beta)}{P(\mathbf{Z}_{-(m,n)}, \mathbf{W}; \alpha, \beta)}$$

- $Z_{(m,n)}$ :  $Z$  hidden variable of the  $n$ -th word token in the  $m$ -th document
- The denominator term does not affect the likelihood:  $\propto P(Z_{(m,n)} = v, \mathbf{Z}_{-(m,n)}, \mathbf{W}; \alpha, \beta)$





# LDA – Gibbs Sampling

## Collapsed Gibbs Sampling

Find the most likely assignment on  $Z$ .

- Gibbs sampling formula

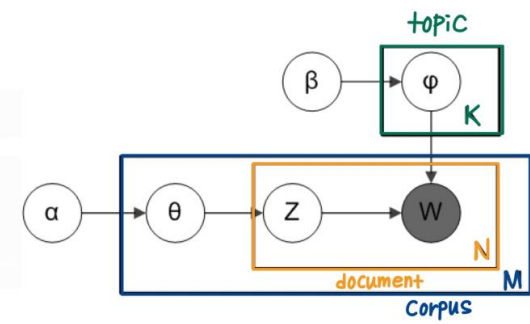
$$P(\mathbf{Z}, \mathbf{W}; \alpha, \beta) = \prod_{j=1}^M \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \frac{\prod_{i=1}^K \Gamma(n_{j,(\cdot)}^i + \alpha_i)}{\Gamma(\sum_{i=1}^K n_{j,(\cdot)}^i + \alpha_i)} \times \prod_{i=1}^K \frac{\Gamma(\sum_{r=1}^V \beta_r)}{\prod_{r=1}^V \Gamma(\beta_r)} \frac{\prod_{r=1}^V \Gamma(n_{(\cdot),r}^i + \beta_r)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^i + \beta_r)}.$$

$$P(Z_{(m,n)} = v \mid \mathbf{Z}_{-(m,n)}, \mathbf{W}; \alpha, \beta)$$

$$\propto P(Z_{(m,n)} = v, \mathbf{Z}_{-(m,n)}, \mathbf{W}; \alpha, \beta)$$

$$\begin{aligned} &= \left( \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \right)^M \prod_{j \neq m} \frac{\prod_{i=1}^K \Gamma(n_{j,(\cdot)}^i + \alpha_i)}{\Gamma(\sum_{i=1}^K n_{j,(\cdot)}^i + \alpha_i)} \left( \frac{\Gamma(\sum_{r=1}^V \beta_r)}{\prod_{r=1}^V \Gamma(\beta_r)} \right)^K \prod_{i=1}^K \left\{ \prod_{r \neq v} \frac{\Gamma(n_{(\cdot),r}^i + \beta_r)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^i + \beta_r)} \right\} \frac{\prod_{i=1}^K \Gamma(n_{m,(\cdot)}^i + \alpha_i)}{\Gamma(\sum_{i=1}^K n_{m,(\cdot)}^i + \alpha_i)} \prod_{i=1}^K \frac{\Gamma(n_{(\cdot),v}^i + \beta_v)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^i + \beta_r)} \\ &\stackrel{\substack{\uparrow \\ j=m \\ r=v}}{\propto} \frac{\prod_{i=1}^K \Gamma(n_{m,(\cdot)}^i + \alpha_i)}{\Gamma(\sum_{i=1}^K n_{m,(\cdot)}^i + \alpha_i)} \prod_{i=1}^K \frac{\Gamma(n_{(\cdot),v}^i + \beta_v)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^i + \beta_r)} \\ &\stackrel{\substack{\uparrow \\ \text{constant term}}}{\propto} \prod_{i=1}^K \Gamma(n_{m,(\cdot)}^i + \alpha_i) \prod_{i=1}^K \frac{\Gamma(n_{(\cdot),v}^i + \beta_v)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^i + \beta_r)}. \end{aligned}$$

$j=m$        $r=v$



# LDA – Gibbs Sampling

## Collapsed Gibbs Sampling

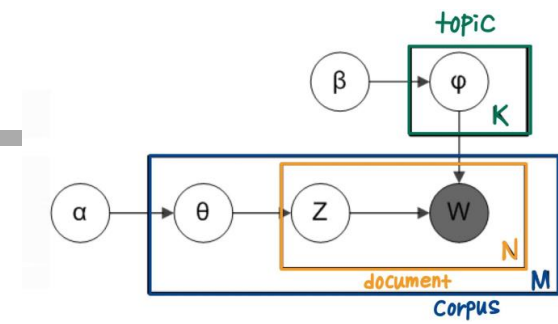
Find the most likely assignment on  $Z$ .

- Gibbs sampling formula

- $n_{j,r}^{i,-(m,n)}$ : the number of word tokens in the  $j$ th document with the same word symbol (the  $r$ th word in the (total) vocabulary) assigned to the  $i$ th topic

- But with  $Z_{(m,n)}$  excluded

- Use property of gamma function:  $\Gamma(x + 1) = x \times \Gamma(x)$



Suppose that current word  $w_i$  is allocated to  $z_i = k$

$$\begin{aligned}
 & \prod_{i=1}^K \Gamma(n_{m,(\cdot)}^i + \alpha_i) \prod_{i=1}^K \frac{\Gamma(n_{(\cdot),v}^i + \beta_v)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^i + \beta_r)} \propto \prod_{i \neq k} \Gamma(n_{m,(\cdot)}^{i,-(m,n)} + \alpha_i) \prod_{i \neq k} \frac{\Gamma(n_{(\cdot),v}^{i,-(m,n)} + \beta_v)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^{i,-(m,n)} + \beta_r)} \frac{\Gamma(n_{m,(\cdot)}^{k,-(m,n)} + \alpha_k + 1)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^{k,-(m,n)} + \beta_r + 1)} \frac{\Gamma(n_{(\cdot),v}^{k,-(m,n)} + \beta_v + 1)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^{k,-(m,n)} + \beta_r + 1)} \\
 & = \prod_{i \neq k} \Gamma(n_{m,(\cdot)}^{i,-(m,n)} + \alpha_i) \prod_{i \neq k} \frac{\Gamma(n_{(\cdot),v}^{i,-(m,n)} + \beta_v)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^{i,-(m,n)} + \beta_r)} \frac{\Gamma(n_{m,(\cdot)}^{k,-(m,n)} + \alpha_k)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^{k,-(m,n)} + \beta_r)} \frac{\Gamma(n_{(\cdot),v}^{k,-(m,n)} + \beta_v)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^{k,-(m,n)} + \beta_r)} \frac{(n_{m,(\cdot)}^{k,-(m,n)} + \alpha_k)}{\sum_{r=1}^V n_{(\cdot),r}^{k,-(m,n)} + \beta_r} \frac{n_{(\cdot),v}^{k,-(m,n)} + \beta_v}{\sum_{r=1}^V n_{(\cdot),r}^{k,-(m,n)} + \beta_r} \\
 & = \prod_i \Gamma(n_{m,(\cdot)}^{i,-(m,n)} + \alpha_i) \prod_i \frac{\Gamma(n_{(\cdot),v}^{i,-(m,n)} + \beta_v)}{\Gamma(\sum_{r=1}^V n_{(\cdot),r}^{i,-(m,n)} + \beta_r)} \frac{(n_{m,(\cdot)}^{k,-(m,n)} + \alpha_k)}{\sum_{r=1}^V n_{(\cdot),r}^{k,-(m,n)} + \beta_r} \frac{n_{(\cdot),v}^{k,-(m,n)} + \beta_v}{\sum_{r=1}^V n_{(\cdot),r}^{k,-(m,n)} + \beta_r} \\
 & \propto (n_{m,(\cdot)}^{k,-(m,n)} + \alpha_k) \frac{n_{(\cdot),v}^{k,-(m,n)} + \beta_v}{\sum_{r=1}^V n_{(\cdot),r}^{k,-(m,n)} + \beta_r}
 \end{aligned}$$