

# Sampling Methods (contd)

CS772A: Probabilistic Machine Learning

Piyush Rai

# Markov Chain Monte Carlo (MCMC)

If the target is a posterior, it will be conditioned on data, i.e.,  $p(\mathbf{z}|\mathbf{x})$

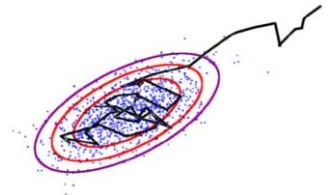
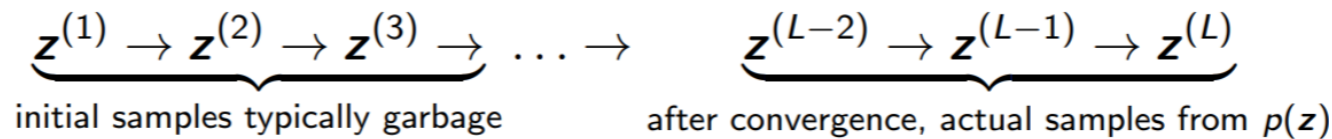
- Goal: Generate samples from some target distribution  $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$

$\mathbf{z}$  usually is high-dim

- Assume we can evaluate  $p(\mathbf{z})$  at least up to a proportionality constant

Means we can at least evaluate  $\tilde{p}(\mathbf{z})$

- MCMC uses a **Markov Chain** which, when converged, starts giving samples from  $p(\mathbf{z})$



- Given current sample  $\mathbf{z}^{(\ell)}$  from the chain, MCMC generates the next sample  $\mathbf{z}^{(\ell+1)}$  as

- Use a **proposal distribution**  $q(\mathbf{z}|\mathbf{z}^{(\ell)})$  to generate a candidate sample  $\mathbf{z}_*$
- **Accept/reject**  $\mathbf{z}_*$  as the next sample based on an **acceptance criterion** (will see later)
- If accepted, set  $\mathbf{z}^{(\ell+1)} = \mathbf{z}_*$ . If rejected, set  $\mathbf{z}^{(\ell+1)} = \mathbf{z}^{(\ell)}$

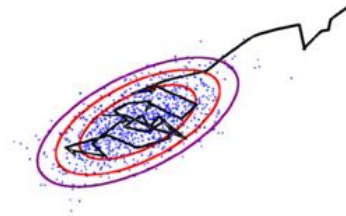
Should also have the same support as  $p(\mathbf{z})$

- Important: The proposal distribution  $q(\mathbf{z}|\mathbf{z}^{(\ell)})$  depends on the previous sample  $\mathbf{z}^{(\ell)}$



# MCMC: The Basic Scheme

3



- The chain run infinitely long (i.e., upon convergence) will give ONE sample from  $p(\mathbf{z})$
- But we usually require **several samples** to approximate  $p(\mathbf{z})$
- This is done as follows
  - Start the chain at an initial  $\mathbf{z}^{(0)}$
  - Using the proposal  $q(\mathbf{z}|\mathbf{z}^{(\ell)})$ , run the chain long enough, say  $T_1$  steps
  - Discard the first  $T_1 - 1$  samples (called “**burn-in**” **samples**) and take last sample  $\mathbf{z}^{(T_1)}$
  - Continue from  $\mathbf{z}^{(T_1)}$  up to  $T_2$  steps, discard intermediate samples, take last sample  $\mathbf{z}^{(T_2)}$ 
    - This discarding (called “**thinning**”) helps ensure that  $\mathbf{z}^{(T_1)}$  and  $\mathbf{z}^{(T_2)}$  are **uncorrelated**
  - Repeat the same for a total of  $S$  times
  - In the end, we now have  $S$  *approximately independent* samples from  $p(\mathbf{z})$
- Note: Good choices for  $T_1$  and  $T_i - T_{i-1}$  (thinning gap) are usually based on heuristics

MCMC is exact in theory but approximate in practice since we can't run the chain for infinitely long in practice

Thus we say that the samples are approximately from the target distribution

Will treat it as our first sample from  $p(\mathbf{z})$

Requirement for Monte Carlo approximation

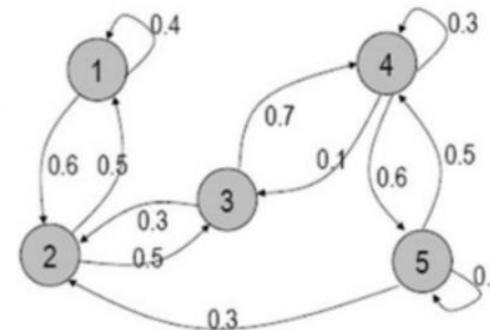


# MCMC: Some Basic Theory

- A first order Markov Chain assumes  $p(\mathbf{z}^{(\ell+1)} | \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(\ell)}) = p(\mathbf{z}^{(\ell+1)} | \mathbf{z}^{(\ell)})$
- A 1st order Markov Chain  $\mathbf{z}^{(0)}, \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(L)}$  is a sequence of r.v.'s and is defined by
  - An **initial state distribution**  $p(\mathbf{z}^{(0)})$
  - A **Transition Function** (TF):  $T_\ell(\mathbf{z}^{(\ell)} \rightarrow \mathbf{z}^{(\ell+1)}) = p(\mathbf{z}^{(\ell+1)} | \mathbf{z}^{(\ell)})$  – the proposal distribution
- TF is a distribution over the values of next state given the value of the current state
- Assuming  $\mathbf{z}$  is discrete with  $K$  possible values, the TF will be  $K \times K$  probability table

Transition probabilities  
can be defined using a  
 $K \times K$  table if  $\mathbf{z}$  is a discrete  
r.v. with  $K$  possible values

	1	2	3	4	5
1	0.4	0.6	0.0	0.0	0.0
2	0.5	0.0	0.5	0.0	0.0
3	0.0	0.3	0.0	0.7	0.0
4	0.0	0.0	0.1	0.3	0.6
5	0.0	0.3	0.0	0.5	0.2



- Homogeneous Markov Chain: The TF is the same for all  $\ell$ , i.e.,  $T_\ell = T$



# MCMC: Some Basic Theory

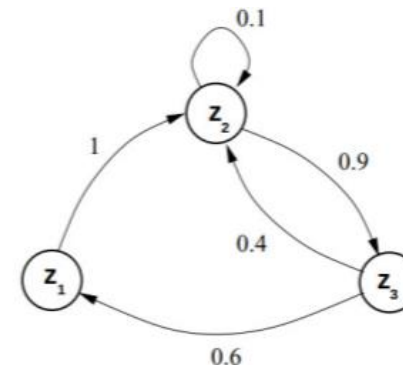
- Consider the following Markov Chain to sample a discrete r.v.  $\mathbf{z}$  with 3 possible values

The initial state distribution for  $\mathbf{z}$

$$p(\mathbf{z}^{(0)}) = p(z_1^{(0)}, z_2^{(0)}, z_3^{(0)}) = [0.5, 0.2, 0.3]$$

Probabilities of the initial state taking each of the 3 possible values

$$T = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0.1 & 0.9 \\ 0.6 & 0.4 & 0 \end{bmatrix}$$



Distribution of  $\mathbf{z}$  after taking the first step

$$p(\mathbf{z}^{(1)}) = p(\mathbf{z}^{(0)}) \times T = [0.2, 0.6, 0.2] \quad (\text{rounded to single digit after decimal})$$

After doing it a few more (say some  $m$ ) times

$$p(\mathbf{z}^{(0)}) \times T^m = [0.2, 0.4, 0.4] \quad (\text{rounded to single digit after decimal})$$

Stationary/Invariant Distribution  
 $p(\mathbf{z})$  of this Markov Chain

$p(\mathbf{z})$  is multinoulli with  $\pi = [0.2, 0.4, 0.4]$

- $p(\mathbf{z})$  being Stationary means no matter what  $p(\mathbf{z}^{(0)})$  is, we will reach  $p(\mathbf{z})$
- Such transition functions are desirable in MCMC



# MCMC: Some Basic Theory

- A Markov Chain with transition function  $T$  has stationary distribution  $p(\mathbf{z})$  if  $T$  satisfies

Known as the Detailed Balance condition

$$p(\mathbf{z})T(\mathbf{z}'|\mathbf{z}) = p(\mathbf{z}')T(\mathbf{z}|\mathbf{z}')$$

Here  $T(b|a)$  denotes the transition probability of going from state  $a$  to state  $b$

- Detailed Balance ensures “reversibility”
- Integrating out (or summing over) detailed balanced condition on both sides w.r.t.  $\mathbf{z}'$

Thus  $p(\mathbf{z})$  is the stationary distribution of this Markov Chain

$$p(\mathbf{z}) = \int p(\mathbf{z}')T(\mathbf{z}|\mathbf{z}')d\mathbf{z}'$$

- Thus a Markov Chain with detailed balance always converges to a stationary distribution
- Detailed balance is sufficient but not necessary condition for having a stationary distr.



# Some MCMC Algorithms



# Metropolis-Hastings (MH) Sampling (1960)

- Suppose we wish to generate samples from a target distribution  $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$
- Assume a suitable proposal distribution  $q(\mathbf{z}|\mathbf{z}^{(\tau)})$ , e.g.,  $\mathcal{N}(\mathbf{z}|\mathbf{z}^{(\tau)}, \sigma^2 \mathbf{I})$
- In each step, draw  $\mathbf{z}^*$  from  $q(\mathbf{z}|\mathbf{z}^{(\tau)})$  and accept  $\mathbf{z}^*$  with probability

$$A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) = \min \left( 1, \frac{\tilde{p}(\mathbf{z}^*) q(\mathbf{z}^{(\tau)}|\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(\tau)}) q(\mathbf{z}^*|\mathbf{z}^{(\tau)})} \right)$$

Favors acceptance of  $\mathbf{z}^*$  if it is more probable than  $\mathbf{z}^{(\tau)}$  (under  $p(\mathbf{z})$ )

Downweight the probability of acceptance of  $\mathbf{z}^*$  if the proposal itself favors its generation (i.e., if  $q(\mathbf{z}^*|\mathbf{z}^{(\tau)})$  is high), and upweight if it unfavors the generation

- Transition function of this Markov Chain

- $T(\mathbf{z}^*|\mathbf{z}^{(\tau)}) = A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) q(\mathbf{z}^*|\mathbf{z}^{(\tau)})$  if state changed
- $T(\mathbf{z}^*|\mathbf{z}^{(\tau)}) = q(\mathbf{z}^{(\tau)}|\mathbf{z}^{(\tau)}) + \sum_{\mathbf{z}^* \neq \mathbf{z}^{(\tau)}} (1 - A(\mathbf{z}^*, \mathbf{z}^{(\tau)})) q(\mathbf{z}^*|\mathbf{z}^{(\tau)})$  otherwise

Can Show that this TF satisfied detailed balance





# The MH Sampling Algorithm

- Initialize  $\mathbf{z}^{(1)}$  randomly
- For  $\ell = 1, 2, \dots, L$ 
  - Sample  $\mathbf{z}^* \sim q(\mathbf{z}^* | \mathbf{z}^{(\ell)})$  and  $u \sim \text{Unif}(0,1)$
  - Compute acceptance probability

$$A(\mathbf{z}^*, \mathbf{z}^{(\ell)}) = \min \left( 1, \frac{\tilde{p}(\mathbf{z}^*)q(\mathbf{z}^{(\ell)} | \mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(\ell)})q(\mathbf{z}^* | \mathbf{z}^{(\ell)})} \right)$$

- If  $A(\mathbf{z}^*, \mathbf{z}^{(\ell)}) > u$

$$\mathbf{z}^{(\ell+1)} = \mathbf{z}^*$$

Meaning accepting  $\mathbf{z}^*$  with probability  $A(\mathbf{z}^*, \mathbf{z}^{(\ell)})$

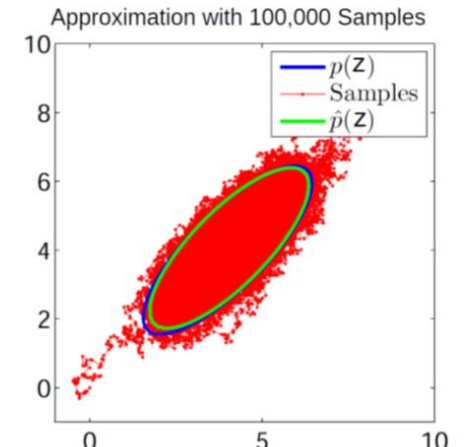
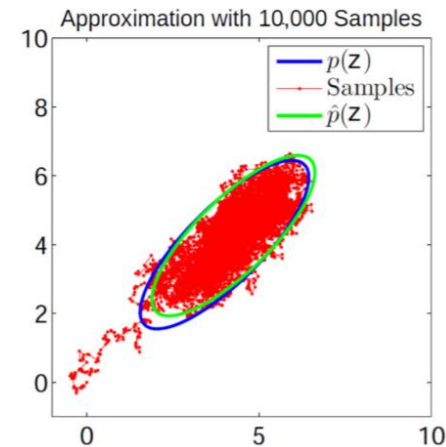
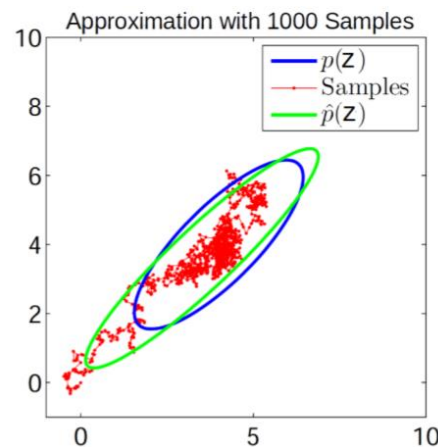
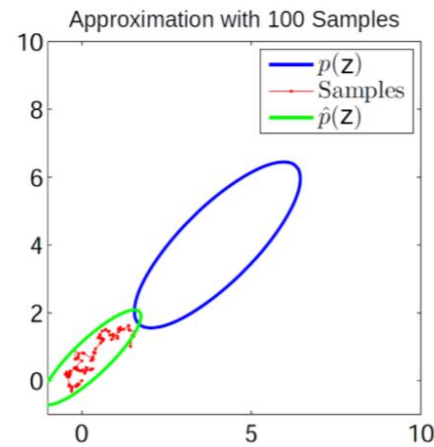
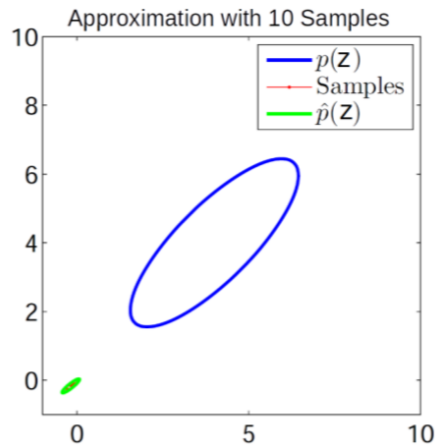
- Else

$$\mathbf{z}^{(\ell+1)} = \mathbf{z}^{(\ell)}$$



# MH Sampling in Action: A Toy Example..

- Target distribution  $p(\mathbf{z}) = \mathcal{N} \left( \begin{bmatrix} 4 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix} \right)$
- Proposal distribution  $q(\mathbf{z}^{(t)} | \mathbf{z}^{(t-1)}) = \mathcal{N} \left( \mathbf{z}^{(t-1)}, \begin{bmatrix} 0.01 & 0 \\ 0 & 0.01 \end{bmatrix} \right)$

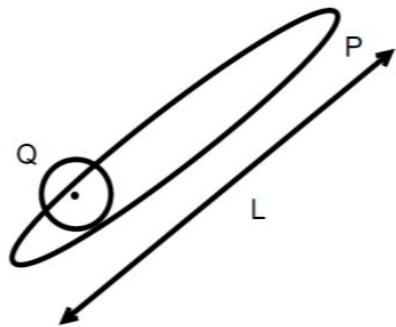


# MH Sampling: Some Comments

- If prop. distrib. is symmetric, we get [Metropolis Sampling](#) algo (Metropolis, 1953) with

$$A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) = \min \left( 1, \frac{\tilde{p}(\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(\tau)})} \right)$$

- Some limitations of MH sampling
  - Can sometimes have very slow convergence (also known as slow “mixing”)



$$Q(\mathbf{z}|\mathbf{z}^{(\tau)}) = \mathcal{N}(\mathbf{z}|\mathbf{z}^{(\tau)}, \sigma^2 \mathbf{I})$$

$\sigma$  large  $\Rightarrow$  many rejections

$\sigma$  small  $\Rightarrow$  slow diffusion

$\sim \left(\frac{L}{\sigma}\right)^2$  iterations required for convergence

- Computing acceptance probability can be expensive\*, e.g., if  $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$  is some target posterior then  $\tilde{p}(\mathbf{z})$  would require computing likelihood on all the data points (expensive)



# Gibbs Sampling (Geman & Geman, 1984)

- Goal: Sample from a joint distribution  $p(\mathbf{z})$  where  $\mathbf{z} = [z_1, z_2, \dots, z_M]$
- Suppose we can't sample from  $p(\mathbf{z})$  but can sample from each conditional  $p(z_i | \mathbf{z}_{-i})$ 
  - In Bayesian models, can be done easily if we have a locally conjugate model
- For Gibbs sampling, the proposal is the conditional distribution  $p(z_i | \mathbf{z}_{-i})$
- Gibbs sampling samples from these conditionals in a cyclic order
- Gibbs sampling is equivalent to MH sampling with acceptance prob. = 1

Hence no need to compute it

$$A(\mathbf{z}^*, \mathbf{z}) = \frac{p(\mathbf{z}^*)q(\mathbf{z}|\mathbf{z}^*)}{p(\mathbf{z})q(\mathbf{z}^*|\mathbf{z})} = \frac{p(z_i^*|\mathbf{z}_{-i}^*)p(\mathbf{z}_{-i}^*)p(z_i|\mathbf{z}_{-i}^*)}{p(z_i|\mathbf{z}_{-i})p(\mathbf{z}_{-i})p(z_i^*|\mathbf{z}_{-i})} = 1$$

where we use the fact that  $\mathbf{z}_{-i}^* = \mathbf{z}_{-i}$

Since only one component is changed at a time



# Gibbs Sampling: Sketch of the Algorithm

- $M$ : Total number of variables,  $T$ : number of Gibbs sampling iterations

1. Initialize  $\{z_i : i = 1, \dots, M\}$  Assuming  $\mathbf{z} = [z_1, z_2, \dots, z_M]$

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  $z_M^{(\tau+1)} \sim p(z_M | z_1^{(\tau+1)}, z_2^{(\tau+1)}, \dots, z_{M-1}^{(\tau+1)})$ .

CP of each component of  $\mathbf{z}$  uses the most recent values (from this or the previous iteration) of all the other components

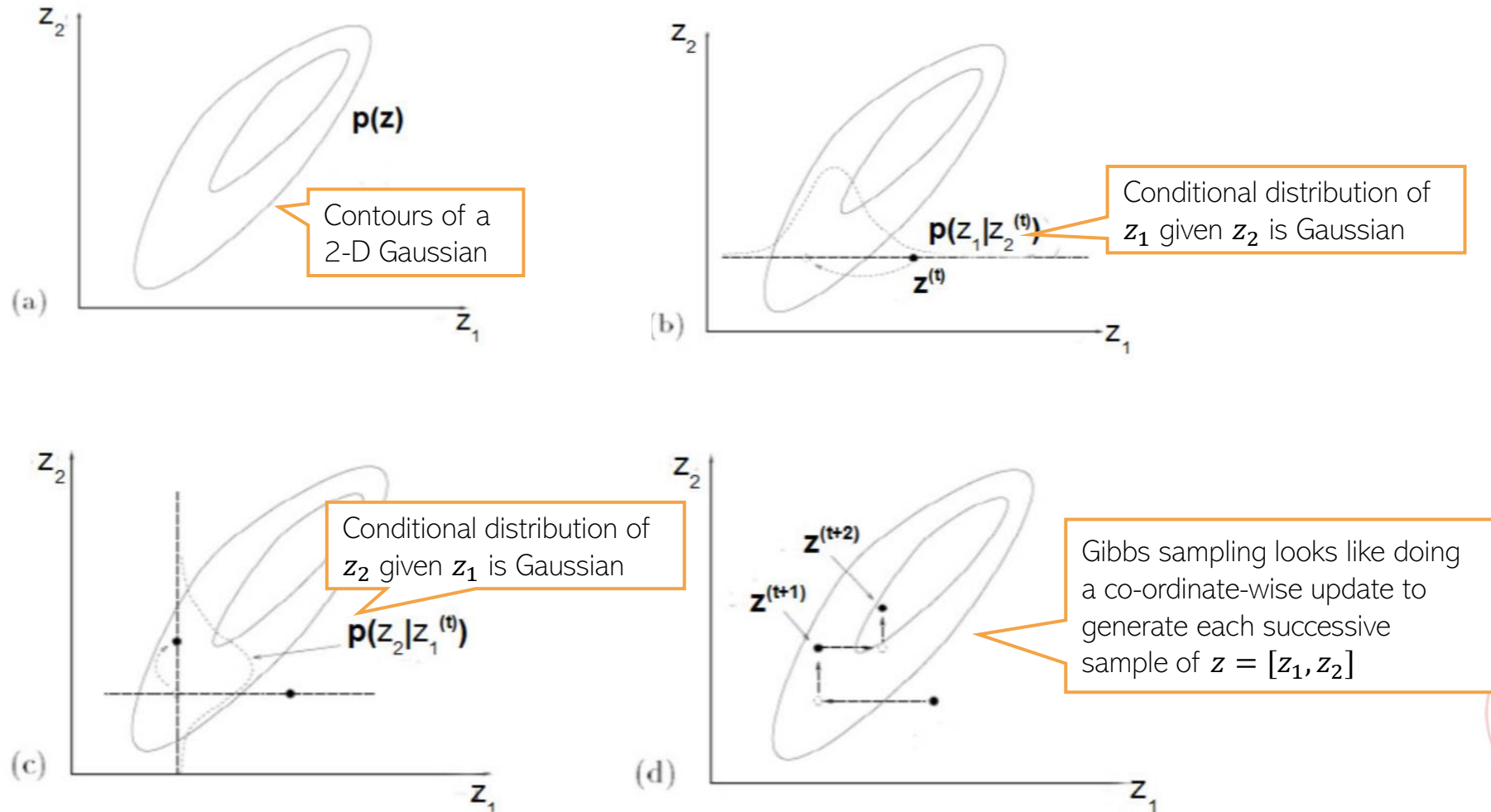
Each iteration will give us one sample  $\mathbf{z}^{(\tau)}$  of  $\mathbf{z} = [z_1, z_2, \dots, z_M]$

- Note: Order of updating the variables usually doesn't matter (but see "Scan Order in Gibbs Sampling: Models in Which it Matters and Bounds on How Much" from NIPS 2016)



# Gibbs Sampling: A Simple Example

- Can sample from a 2-D Gaussian using 1-D Gaussians



# Gibbs Sampling: Another Simple Example

- Bayesian linear regression:  $p(y_n | \mathbf{x}_n, \mathbf{w}, \beta) = \mathcal{N}(y_n | \mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$ ,  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda^{-1} I)$ ,  $p(\lambda) = \text{Gamma}(\lambda | a, b)$ ,  $p(\beta) = \text{Gamma}(\beta | c, d)$ . Gibbs sampler for  $p(\mathbf{w}, \lambda, \beta | \mathbf{X}, \mathbf{y})$  will be
- Initialize  $\lambda, \beta$  as  $\lambda^{(0)}, \beta^{(0)}$ . For iteration  $t = 1, 2, \dots, T$

- Generate a random sample of  $\mathbf{w}$  by sampling from its CP as

$$\mathbf{w}^{(t)} \sim \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)}) \quad \text{where}$$

$$\boldsymbol{\Sigma}^{(t-1)} = (\beta^{(t-1)} \mathbf{X}^\top \mathbf{X} + \lambda^{(t-1)})^{-1}$$

$$\boldsymbol{\mu}^{(t-1)} = \left( \mathbf{X}^\top \mathbf{X} + \frac{\lambda^{(t-1)}}{\beta^{(t-1)}} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

- Generate a random sample of  $\lambda$  by sampling from its CP as

$$\lambda^{(t)} \sim \text{Gamma} \left( \lambda | a + \frac{D}{2}, b + \frac{\mathbf{w}^{(t)\top} \mathbf{w}^{(t)}}{2} \right)$$

- Generate a random sample of  $\beta$  by sampling from its CP as

$$\beta^{(t)} \sim \text{Gamma} \left( \beta | c + \frac{N}{2}, d + \frac{\|\mathbf{y} - \mathbf{X} \mathbf{w}^{(t)}\|^2}{2} \right)$$

Note: Assuming these are post-burnin samples and thinning (if any) is also considered

- The posterior's approximation is the set of collected samples  $\{\mathbf{w}^{(t)}, \lambda^{(t)}, \beta^{(t)}\}_{t=1}^T$

