

# Chapter 17

## Monte Carlo Methods

## Monte Carlo Sampling

---

- To approximate sums or integrals (which are costly to evaluate or intractable) by drawing samples

$$s = \sum_{\mathbf{x}} p(\mathbf{x}) f(\mathbf{x}) \text{ or } s = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$

- **Idea:** To view the sum/integral as an expectation under some distribution and to approximate it by an *average*

$$s = E_p[f(\mathbf{x})] \approx \hat{s}_n = \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^{(i)})$$

where

$$\mathbf{x}^{(i)} \sim p(\mathbf{x})$$

- It is easy to verify that the estimator  $\hat{s}_n$  is unbiased

$$E[\hat{s}_n] = E_p[f(\mathbf{x})] = s$$

- If the samples  $\mathbf{x}^{(i)}$  are independently and identically distributed (i.i.d.),

$$\text{Var}[\hat{s}_n] = \frac{\text{Var}[f(\mathbf{x})]}{n}$$

$$\hat{s}_n \sim \mathcal{N}(s, \text{Var}[\hat{s}_n]) \quad (\text{C.L.T.})$$

## Importance Sampling

- To approximate the expectation based on a **proposal distribution**  $q(\mathbf{x})$  that is easier to draw samples from than  $p(\mathbf{x})$

$$s = \sum_{\mathbf{x}} p(\mathbf{x}) f(\mathbf{x}) = \sum_{\mathbf{x}} q(\mathbf{x}) \frac{p(\mathbf{x}) f(\mathbf{x})}{q(\mathbf{x})}$$

- Importance sampling estimator  $\hat{s}_q$

$$\hat{s}_q = \frac{1}{n} \sum_{i=1}^n \frac{p(\mathbf{x}^{(i)}) f(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} = \frac{1}{n} \sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})$$

where

$$\mathbf{x}^{(i)} \sim q(\mathbf{x})$$

- $p(\mathbf{x}^{(i)})/q(\mathbf{x}^{(i)})$  are known as *importance weights*

- It is readily seen that  $\hat{s}_q$  is unbiased irrespective of the choice of  $q(\mathbf{x})$

$$E_q[\hat{s}_q] = E_q\left[\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right] = E_p[f(\mathbf{x})] = s$$

- The variance of  $\hat{s}_q$  is however highly sensitive to the choice of  $q(\mathbf{x})$

$$\text{Var}[\hat{s}_q] = \text{Var}\left[\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right]/n$$

## Biased Importance Sampling

---

- Oftentimes  $p(\mathbf{x})$  can only be evaluated up to a normalization constant

$$p(\mathbf{x}) = \frac{\tilde{p}(\mathbf{x})}{Z_p}$$

That is,  $\tilde{p}(\mathbf{x})$  is easy to evaluate and  $Z_p$  is unknown (or intractable)

- We may also wish to use a  $q(\mathbf{x})$  with the same property

$$q(\mathbf{x}) = \frac{\tilde{q}(\mathbf{x})}{Z_q}$$

- The importance sampling estimator is then given by

$$\begin{aligned}\hat{s}_q &= \frac{1}{n} \sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)}) \\ &= \frac{Z_q}{Z_p} \frac{1}{n} \sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})\end{aligned}$$

$$= \frac{Z_q}{Z_p} \frac{1}{n} \sum_{i=1}^n \tilde{r}_i f(\mathbf{x}^{(i)})$$

where

$$\tilde{r}_i = \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} \text{ and } \mathbf{x}^{(i)} \sim q(\mathbf{x})$$

- The same set of data  $\mathbf{x}^{(i)}$  can be used to approximate the ratio  $Z_p/Z_q$

$$\begin{aligned} \frac{Z_p}{Z_q} &= \frac{\sum_{\mathbf{x}} \tilde{p}(\mathbf{x})}{Z_q} \\ &= \sum_{\mathbf{x}} \tilde{p}(\mathbf{x}) \frac{1}{Z_q} \\ &= \sum_{\mathbf{x}} \tilde{p}(\mathbf{x}) \frac{q(\mathbf{x})}{\tilde{q}(\mathbf{x})} \\ &= \sum_{\mathbf{x}} \frac{\tilde{p}(\mathbf{x})}{\tilde{q}(\mathbf{x})} q(\mathbf{x}) \end{aligned}$$

$$\begin{aligned} &\simeq \frac{1}{n} \sum_i \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} \\ &= \frac{1}{n} \sum_i \tilde{r}_i \end{aligned}$$

- We then arrive at a *biased importance estimator*

$$\hat{s}_{BIS} = \frac{\sum_{i=1}^n \tilde{r}_i f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \tilde{r}_i} = \sum_{i=1}^n \tilde{w}_i f(\mathbf{x}^{(i)})$$

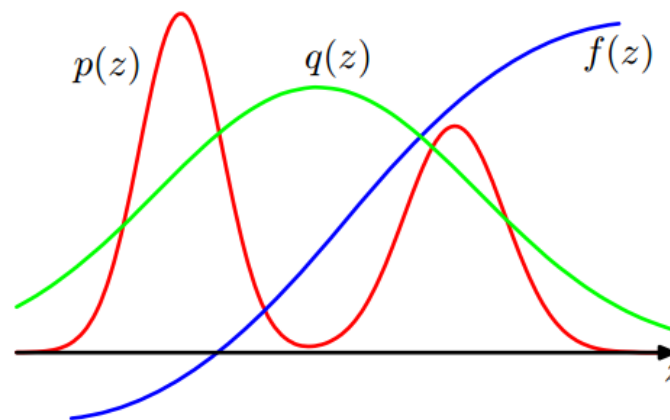
where

$$\tilde{w}_i = \frac{\tilde{r}_i}{\sum_{i=1}^n \tilde{r}_i}$$

- $\hat{s}_{BIS}$  is asymptotically unbiased; that is, as  $n \rightarrow \infty$ ,  $E[\hat{s}_{BIS}] = s$



- The success of importance sampling depends crucially on how well  $q(\mathbf{x})$  matches the desired distribution  $p(\mathbf{x})$
- When  $p(\mathbf{x})f(\mathbf{x})$  is strongly varying and has its mass concentrated over small regions of  $\mathbf{x}$  space, most samples collected may be useless since they contribute little to the final estimate due to the fact  $q(\mathbf{x}^{(i)}) \gg p(\mathbf{x}^{(i)})|f(\mathbf{x}^{(i)})|$
- As such, underestimation of  $E_p[f(\mathbf{x})]$  is typical, especially when  $\mathbf{x}$  is high dimensional



## Markov Chain Monte Carlo Methods

---

- Methods that involve drawing samples from Markov chains to perform Monte Carlo estimation
- Drawing samples from a Markov Chain
  1. Start with an initial state  $\mathbf{x}^{(1)}$
  2. Sample repeatedly from transition distributions  $p(\mathbf{x}^{(\tau+1)}|\mathbf{x}^{(\tau)})$

$$\text{Sample } \mathbf{x}^{(\tau+1)} \sim p(\mathbf{x}^{(\tau+1)}|\mathbf{x}^{(\tau)}), \tau = 1, \dots, t-1$$

- Given a desired distribution  $p^*(\mathbf{x})$ , we choose transition distributions such that  $\mathbf{x}^{(t)}$  eventually becomes a fair sample of  $p^*(\mathbf{x})$

## First-Order Markov Chains

- A sequence of discrete-valued random variables  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}$  with the conditional independence property

$$p(\mathbf{x}^{(m+1)} | \mathbf{x}^{(m)}, \dots, \mathbf{x}^{(1)}) = p(\mathbf{x}^{(m+1)} | \mathbf{x}^{(m)}),$$

for  $m \in \{1, \dots, M - 1\}$

- The joint distribution of  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}$  is characterized by  $p(\mathbf{x}^{(1)})$  together with the transition probabilities  $p(\mathbf{x}^{(m+1)} | \mathbf{x}^{(m)})$

$$p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}) = p(\mathbf{x}^{(1)}) \prod_{i=1}^{M-1} p(\mathbf{x}^{(i+1)} | \mathbf{x}^{(i)})$$



- The marginal distribution  $p(\mathbf{x}^{(m+1)})$  can be expressed as

$$p(\mathbf{x}^{(m+1)}) = \sum_{\mathbf{x}^{(m)}} p(\mathbf{x}^{(m+1)} | \mathbf{x}^{(m)}) p(\mathbf{x}^{(m)})$$

- In matrix form, we have

$$\mathbf{v}^{(m+1)} = \mathbf{A}^{(m)} \mathbf{v}^{(m)}$$

where

$$v_i^{(m+1)} = p(\mathbf{x}^{(m+1)} = \mathbf{s}_i), \quad \text{Prob. of } \mathbf{x}^{(m+1)} \text{ in state } \mathbf{s}_i$$

$$v_j^{(m)} = p(\mathbf{x}^{(m)} = \mathbf{s}_j), \quad \text{Prob. of } \mathbf{x}^{(m)} \text{ in state } \mathbf{s}_j$$

$$A_{i,j}^{(m)} = p(\mathbf{x}^{(m+1)} = \mathbf{s}_i | \mathbf{x}^{(m)} = \mathbf{s}_j), \quad \text{Transition probabilities}$$

- A Markov chain is said to be **homogeneous** if the transition probability  $p(\mathbf{x}^{(m+1)} | \mathbf{x}^{(m)})$  does not depend on  $m$

- In this case, we see that  $\mathbf{A}^{(m)} = \mathbf{A}$  is a constant matrix and that over time, all the eigenvalues are exponentiated

$$\mathbf{v}^{(t)} = \mathbf{A}^{t-1} \mathbf{v}^{(1)} = \mathbf{U} \mathbf{\Lambda}^{t-1} \mathbf{U}^{-1} \mathbf{v}^{(1)}$$

- Under some conditions (e.g. non-zero transition probabilities),  $\mathbf{A}$  has only one eigenvector  $\mathbf{v}$  with the largest eigenvalue 1
- $\mathbf{v}^{(t)}$  eventually converges to that eigenvector  $\mathbf{v}$ , which denotes the **equilibrium distribution**, regardless of the choice of initial state  $\mathbf{v}^{(1)}$

$$\mathbf{A}\mathbf{v} = \mathbf{v}$$

- We hope that by choosing transition probabilities correctly,  $\mathbf{v}$  will be equal to the distribution we wish to sample from

- Running the Markov chain until it reaches its equilibrium is called **burning in** and the time required is called the **mixing time**
- Unfortunately, we only know that the chain will converge under some mild conditions, but not how much time it will take
- Most properties of discrete-valued Markov chains as presented here can carry over to the continuous-valued case

## Gibbs Sampling

- To build a Markov chain that samples from a distribution  $p_{\text{model}}(\mathbf{x})$

$$p_{\text{model}}(\mathbf{x}) = p_{\text{model}}(x_1, x_2, \dots, x_M)$$

- Procedure

1. Start with an initial state  $x_i^{(1)}, i = 1, 2, \dots, M$

2. For  $\tau = 1, \dots, t - 1$

- Sample  $x_1^{(\tau+1)} \sim p_{\text{model}}(x_1 | x_2^{(\tau)}, x_3^{(\tau)}, \dots, x_M^{(\tau)})$
- Sample  $x_2^{(\tau+1)} \sim p_{\text{model}}(x_2 | x_1^{(\tau+1)}, x_3^{(\tau)}, \dots, x_M^{(\tau)})$
- $\vdots$
- Sample  $x_j^{(\tau+1)} \sim p_{\text{model}}(x_j | x_1^{(\tau+1)}, \dots, x_{j-1}^{(\tau+1)}, x_{j+1}^{(\tau)}, \dots, x_M^{(\tau)})$
- $\vdots$
- Sample  $x_M^{(\tau+1)} \sim p_{\text{model}}(x_M | x_1^{(\tau+1)}, x_2^{(\tau+1)}, \dots, x_{M-1}^{(\tau+1)})$

- In words, each step replaces one variable  $x_i$  by drawing a sample from the distribution  $p_{\text{model}}(x_i | \mathbf{x}_{-i})$  of  $x_i$  conditioned on the values of the remaining variables  $\mathbf{x}_{-i}$
- This procedure eventually yields samples of  $p_{\text{model}}(\mathbf{x})$  because
  - The resulting Markov chain will converge to an equilibrium distribution, if none of the transition probabilities is zero anywhere
  - $p_{\text{model}}(\mathbf{x})$  is **invariant** w.r.t. this Markov chain
- A distribution  $p^*(\mathbf{x})$  is said to be invariant w.r.t. a Markov chain if each step in the chain leaves that distribution invariant, i.e.

$$p(\mathbf{x}') = \sum_{\mathbf{x}} p(\mathbf{x}' | \mathbf{x}) p^*(\mathbf{x}) = p^*(\mathbf{x}')$$



- In the present case, we have

$$\mathbf{x} = (x_i^{old}, \mathbf{x}_{-i}^{old}) \sim p_{\text{model}}(\mathbf{x})$$

$$\mathbf{x}' = (x_i^{new}, \mathbf{x}_{-i}^{old}) \text{ with } x_i^{new} \sim p_{\text{model}}(x_i | \mathbf{x}_{-i}^{old})$$

- It can be shown that  $p(\mathbf{x}') = p_{\text{model}}(\mathbf{x}')$ ; that is,  $p_{\text{model}}(\mathbf{x})$  is invariant

$$\begin{aligned} p(\mathbf{x}') &= p(x_i^{new}, \mathbf{x}_{-i}^{old}) \\ &= p(\mathbf{x}_{-i}^{old}) p(x_i^{new} | \mathbf{x}_{-i}^{old}) \\ &= p_{\text{model}}(\mathbf{x}_{-i}^{old}) p_{\text{model}}(x_i^{new} | \mathbf{x}_{-i}^{old}) \\ &= p_{\text{model}}(x_i^{new}, \mathbf{x}_{-i}^{old}) \\ &= p_{\text{model}}(\mathbf{x}') \end{aligned}$$

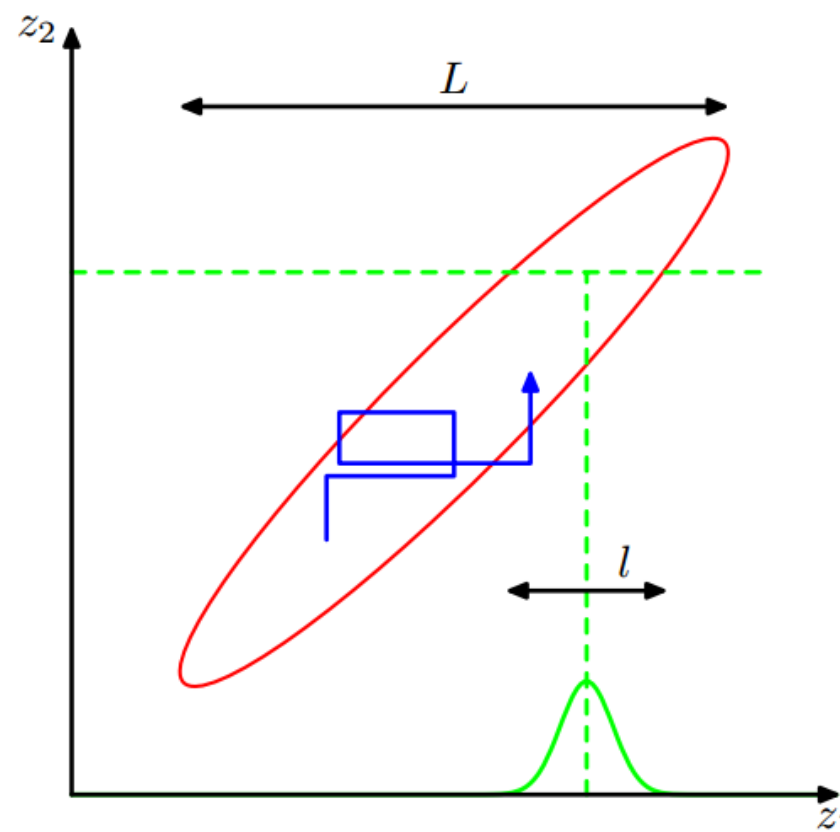
- **Block Gibbs sampling:** In some cases, it is possible to sample many variables simultaneously; for example, in RBM,  $p(\mathbf{h}|\mathbf{v})$  and  $p(\mathbf{v}|\mathbf{h})$  are factorial, suggesting that the elements of  $\mathbf{h}$  and of  $\mathbf{v}$  can be sampled simultaneously

## Challenges

---

- Successive samples are preferably independent and different regions in  $x$  space should be visited proportional to their probability
- In reality, successive samples are highly correlated even though they have identical distributions
- Independent samples may be obtained by retaining every  $M$  samples for sufficiently large  $M$ , or by running multiple chains in parallel

- Moreover, Gibbs sampling may mix slowly when the variables of  $p_{\text{model}}(\mathbf{x})$  are highly correlated



Sampling a correlated Gaussian of two variables

- Mixing between modes may be difficult if they are widely separated by regions of low probability
  - Toy problem: Consider the following energy model

$$\tilde{p}(a, b) = \exp(-E(a, b)), \quad a, b \in \{-1, 1\}$$

where

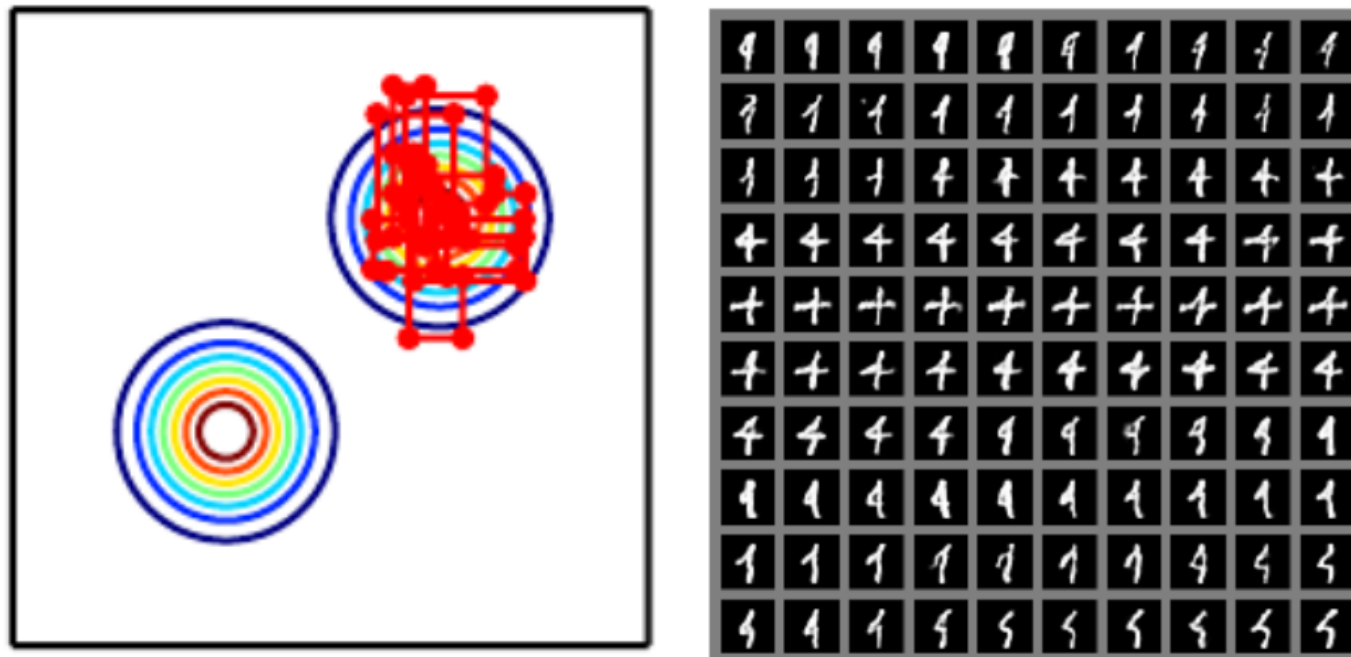
$$E(a, b) = -wab$$

- It is seen that

$$p(b = 1|a = 1) = \sigma(w)$$

- When  $w$  is extremely large, Gibbs sampling will only rarely flip the signs of  $a, b$  even if  $p(b = 1, a = 1) = p(b = -1, a = -1)$

— More examples:



## Confronting The Partition Function

- Many undirected graphical models are defined by an unnormalized distribution  $\tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$  with an intractable partition function  $Z(\boldsymbol{\theta})$

$$p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) = \frac{\tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})}{Z(\boldsymbol{\theta})}$$

where

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) \text{ or } Z(\boldsymbol{\theta}) = \int_{\mathbf{x}} \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}$$

- For training, we maximize the log-likelihood w.r.t. training data

$$E_{\mathbf{x} \sim p_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) = E_{\mathbf{x} \sim p_{\text{data}}} \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) - \log Z(\boldsymbol{\theta})$$

through gradient descent

$$\nabla_{\boldsymbol{\theta}} E_{\mathbf{x} \sim p_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) = \underbrace{E_{\mathbf{x} \sim p_{\text{data}}} \nabla_{\boldsymbol{\theta}} \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})}_{\text{Positive phase}} - \underbrace{\nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta})}_{\text{Negative phase}}$$

- For discrete-valued  $\mathbf{x}$ , the gradient of  $\log Z$  can be evaluated as

$$\nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta}) = \frac{\nabla_{\boldsymbol{\theta}} Z(\boldsymbol{\theta})}{Z(\boldsymbol{\theta})} = \frac{\nabla_{\boldsymbol{\theta}} \sum_{\mathbf{x}} \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})}{Z(\boldsymbol{\theta})} = \frac{\sum_{\mathbf{x}} \nabla_{\boldsymbol{\theta}} \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})}{Z(\boldsymbol{\theta})}$$

- Additionally, if  $\tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) > 0$  for all  $\mathbf{x}$  (e.g. energy-based models),

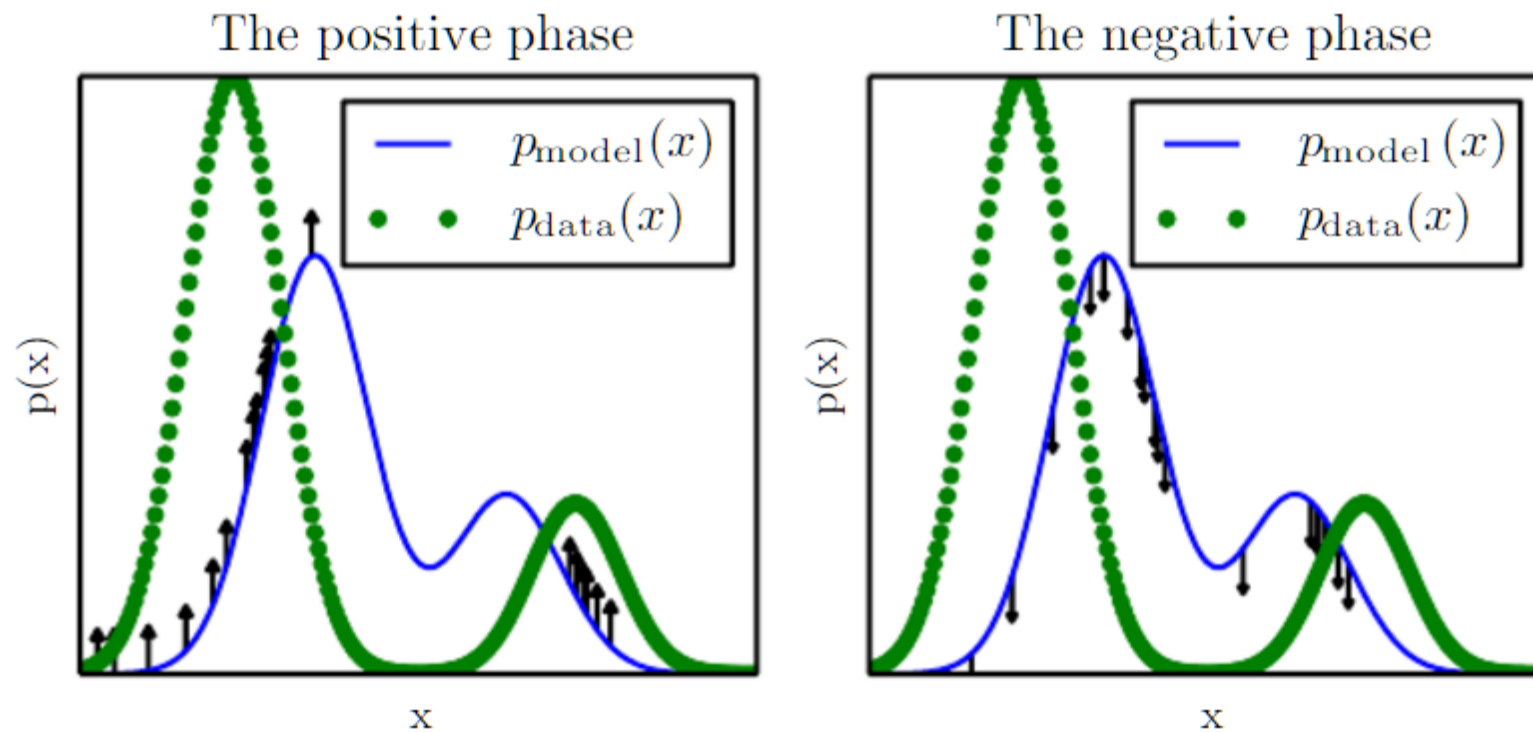
$$\begin{aligned} \frac{\sum_{\mathbf{x}} \nabla_{\boldsymbol{\theta}} \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})}{Z(\boldsymbol{\theta})} &= \frac{\sum_{\mathbf{x}} \nabla_{\boldsymbol{\theta}} \exp(\log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}))}{Z(\boldsymbol{\theta})} \\ &= \frac{\sum_{\mathbf{x}} \exp(\log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}} \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})}{Z(\boldsymbol{\theta})} \\ &= \frac{\sum_{\mathbf{x}} \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})}{Z(\boldsymbol{\theta})} \\ &= \sum_{\mathbf{x}} p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) \\ &= E_{\mathbf{x} \sim p_{\text{model}}} \nabla_{\boldsymbol{\theta}} \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) \end{aligned}$$



- To summarize, we see that

$$\begin{aligned} & \nabla_{\boldsymbol{\theta}} E_{\mathbf{x} \sim p_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) \\ &= E_{\mathbf{x} \sim p_{\text{data}}} \nabla_{\boldsymbol{\theta}} \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) - E_{\mathbf{x} \sim p_{\text{model}}} \nabla_{\boldsymbol{\theta}} \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) \end{aligned}$$

- In the **positive phase**, we increase the log-likelihood by increasing  $\log \tilde{p}(\mathbf{x}; \boldsymbol{\theta})$  with  $\mathbf{x}$  drawn from **training data**  $p_{\text{data}}(\mathbf{x})$
- In the **negative phase**, we increase the log-likelihood by decreasing the partition function  $Z(\boldsymbol{\theta})$ , or equivalently, by decreasing  $\log \tilde{p}(\mathbf{x}; \boldsymbol{\theta})$  with  $\mathbf{x}$  drawn from the **model distribution**  $p_{\text{model}}(\mathbf{x})$
- When  $p_{\text{model}}(\mathbf{x}) = p_{\text{data}}(\mathbf{x})$ , there is no longer gradient



## Contrastive Divergence and Its Variants

---

- To compute the gradient of the negative phase with Gibbs sampling

$$E_{\mathbf{x} \sim p_{\text{model}}} \nabla_{\boldsymbol{\theta}} \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$$

- There are different strategies for initializing the Markov chains
  - Contrastive divergence (CD) – from training data
  - Persistent contrastive divergence (PCD) – from previous step
  - (Study by yourself)

- Example: Contrastive Divergence (CD)

**while** not converged **do**

Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set.

$\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$ .

**for**  $i = 1$  to  $m$  **do**

$\tilde{\mathbf{x}}^{(i)} \leftarrow \mathbf{x}^{(i)}$ .

**end for**

**for**  $i = 1$  to  $k$  **do**

**for**  $j = 1$  to  $m$  **do**

$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)})$ .

**end for**

**end for**

$\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta})$ .

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}$ .

**end while**

## Review

---

- Why sampling?
- Importance sampling
- Gibbs sampling
- Issues with mixing of MCMC methods
- MCMC approach to learning with intractable partition functions