

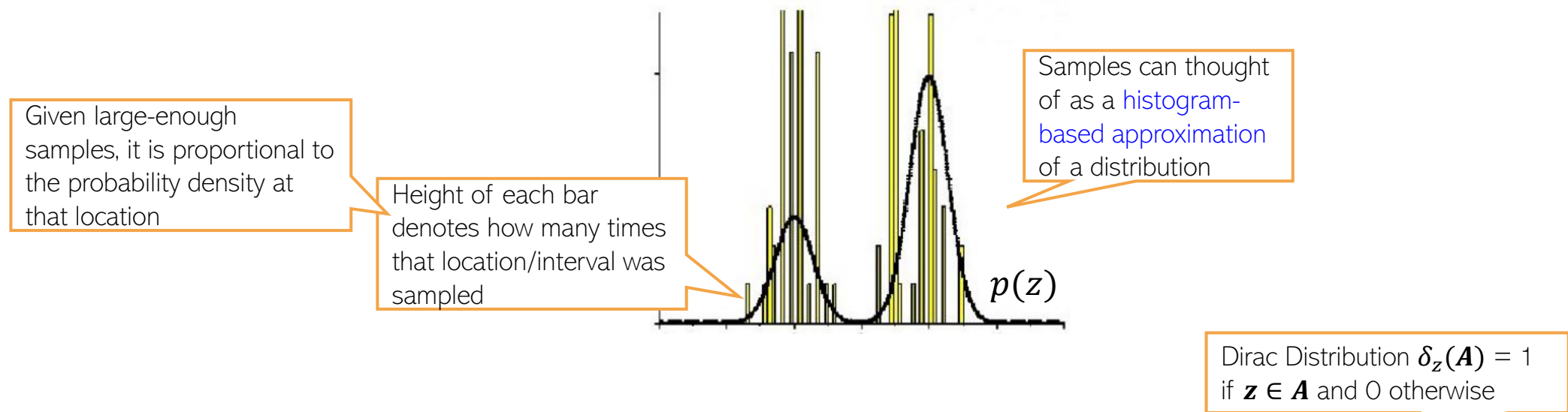
# Approximating Distributions via Sampling

CS772A: Probabilistic Machine Learning

Piyush Rai

# Approximating a Prob. Distribution using Samples <sup>2</sup>

- A distribution can be approximated using a set of **randomly drawn samples** from it



- Using the  $S$  samples  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(S)}$ , our approx.  $p(\mathbf{Z}) \approx \frac{1}{S} \sum_{s=1}^S \delta_{\mathbf{z}^{(s)}}(\mathbf{Z})$
- Usually straightforward to generate samples if it is a simple/standard distribution
- **The interesting bit:** Even if the distribution is “difficult” (e.g., an intractable posterior), it is often possible to generate random samples from such a distribution, as we will see.



# Using samples to compute expectations/predictions

- Given  $S$  i.i.d. samples  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(S)}$  from a distribution  $p(\mathbf{Z})$

$$\mathbb{E}_{p(\mathbf{Z})}[f(\mathbf{Z})] = \int f(\mathbf{Z})p(\mathbf{Z})d\mathbf{Z} \approx \frac{1}{S} \sum_{s=1}^S f(\mathbf{z}^{(s)})$$

- For Bayesian lin. reg., assuming  $\mathbf{w}, \beta, \lambda$  to be unknown, PPD approximation will be

$$\int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) d\mathbf{w} d\beta d\lambda \approx \frac{1}{S} \sum_{s=1}^S p(y_* | \mathbf{x}_*, \mathbf{w}^{(s)}, \beta^{(s)})$$

Joint posterior over all unknowns

Thus, in this case, the PPD is a sum of  $S$  Gaussians

Can also think of it as an **ensemble** consisting of  $S$  members

Sampling based approximation of PPD

Mean:  $\mathbb{E}[\mathbf{w}^T \mathbf{x}_*] \approx \frac{1}{S} \sum_{s=1}^S \mathbf{w}^{(s)T} \mathbf{x}_*$

Variance: Exercise! Use definition of variance and use Monte-Carlo approximation

- Sampling based approx. for PPD of other models can also be obtained likewise



# Sampling: Some Basic Methods

$$p(z) = q(x) \left| \frac{\partial x}{\partial z} \right| \quad 4$$

Determinant of Jacobian

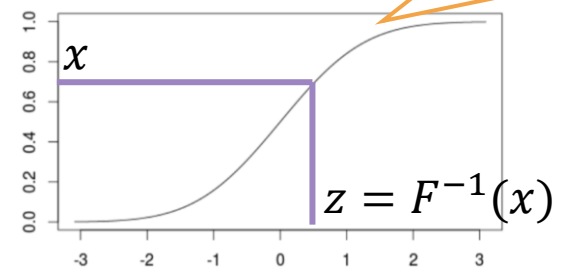
- Most of these basic methods are based on the idea of transformation
  - Generate a random sample  $x$  from a distribution  $q(x)$  which is easy to sample from
  - Apply a transformation on  $x$  to make it random sample  $z$  from a complex distr  $p(z)$

$F(z)$ : CDF of  $p(z)$

- Some popular examples of transformation methods

- Inverse CDF method

$$x \sim \text{Unif}(0, 1) \Rightarrow z = \text{Inv-CDF}_{p(z)}(x) \sim p(z)$$



- Reparametrization method

$$x \sim \mathcal{N}(0, 1) \Rightarrow z = \mu + \sigma x \sim \mathcal{N}(\mu, \sigma^2)$$

- Box-Mueller method: Given  $(x_1, x_2)$  from  $\text{Unif}(0, 1)$ , generate  $(z_1, z_2)$  from  $\mathcal{N}(0, \mathbf{I}_2)$

$$z_1 = \sqrt{-2 \ln x_1} \cos(2\pi x_2), \quad z_2 = \sqrt{-2 \ln x_1} \sin(2\pi x_2)$$

- Transformation Methods are simple but have limitations

- Mostly limited to standard distributions and/or distributions with very few variables

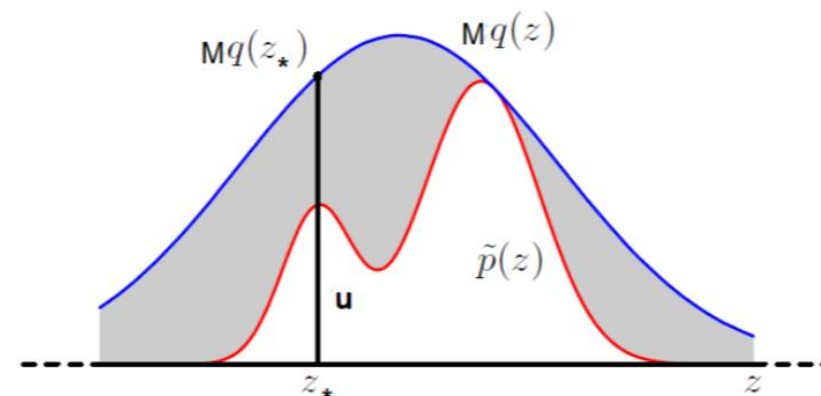


# Rejection Sampling

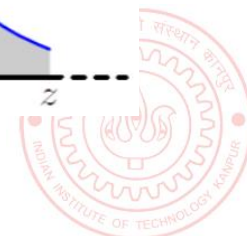
- Goal: Generate a random sample from a distribution of the form  $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$ , assuming
  - We can only evaluate the value of numerator  $\tilde{p}(\mathbf{z})$  for any  $\mathbf{z}$
  - The denominator (normalization constant)  $Z_p$  is intractable and we don't know its value
- Assume a Should have the same support as  $p(\mathbf{z})$  proposal distribution  $q(\mathbf{z})$  we can generate samples from, and

$$Mq(\mathbf{z}) \geq \tilde{p}(\mathbf{z}) \quad \forall \mathbf{z} \quad (\text{where } M > 0 \text{ is some const.})$$

- Rejection Sampling then works as follows
  - Sample a random variable  $\mathbf{z}_*$  from  $q(\mathbf{z})$
  - Sampling a uniform r.v.  $u \sim \text{Unif}[0, Mq(\mathbf{z}_*)]$
  - If  $u \leq \tilde{p}(\mathbf{z}_*)$  then accept  $\mathbf{z}_*$ , otherwise reject it



- All accepted  $\mathbf{z}_*$ 's will be random samples from  $p(\mathbf{z})$ . Proof on next slide



# Rejection Sampling

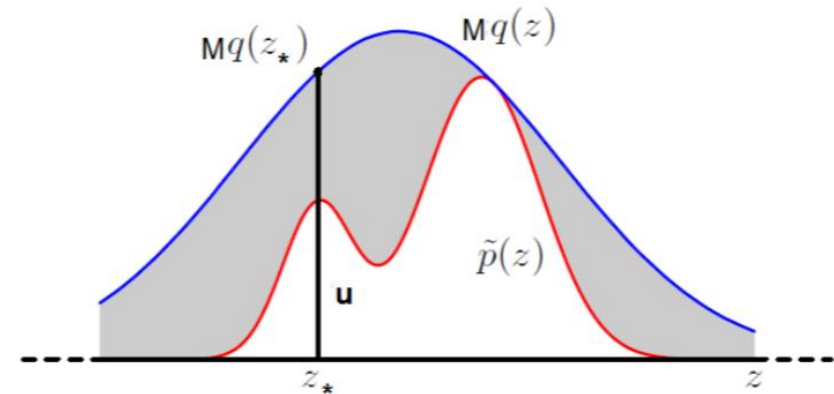
- Why  $z \sim q(z)$  + accept/reject rule is equivalent to  $z \sim p(z)$ ?
- Let's look at the pdf of the  $z$ 's that were accepted, i.e.,  $p(z|\text{accept})$

$$p(\text{accept}|z) = \int_0^{\tilde{p}(z)} \frac{1}{Mq(z)} du = \frac{\tilde{p}(z)}{Mq(z)}$$

$$p(z, \text{accept}) = q(z)p(\text{accept}|z) = \frac{\tilde{p}(z)}{M}$$

$$p(\text{accept}) = \int \frac{\tilde{p}(z)}{M} dz = \frac{Z_p}{M}$$

$$p(z|\text{accept}) = \frac{p(z, \text{accept})}{p(\text{accept})} = \frac{\tilde{p}(z)}{Z_p} = p(z)$$



# Computing Expectations via Monte Carlo Sampling<sup>7</sup>

- Often we are interested in computing expectations of the form

$$\mathbb{E}[f] = \int f(z)p(z)dz$$

where  $f(z)$  is some function of the random variable  $z \sim p(z)$

- A simple approx. scheme to compute the above expectation: [Monte Carlo integration](#)

- Generate  $L$  independent samples from  $p(z)$ :  $\{z^{(\ell)}\}_{\ell=1}^L \sim p(z)$  Assuming we know how to sample from  $p(z)$
- Approximate the expectation by the following empirical average

$$\mathbb{E}[f] \approx \hat{f} = \frac{1}{L} \sum_{\ell=1}^L f(z^{(\ell)})$$

- Since the samples are independent of each other, we can show the following (exercise)

Unbiased expectation

$$\mathbb{E}[\hat{f}] = \mathbb{E}[f]$$

$$\text{and } \text{var}[\hat{f}] = \frac{1}{L} \text{var}[f] = \frac{1}{L} \mathbb{E}[(f - \mathbb{E}[f])^2]$$

Variance in our estimate decreases as  $L$  increases

# Computing Expectations via Importance Sampling <sup>8</sup>

- How to compute Monte Carlo expec. if we don't know how to sample from  $p(\mathbf{z})$ ?
- One way is to use transformation methods or rejection sampling to generate samples
- Another way is to use **Importance Sampling** (assuming  $p(\mathbf{z})$  can be evaluated at least)
  - Generate  $L$  indep samples from a **proposal**  $q(\mathbf{z})$  we know how sample from:  $\{\mathbf{z}^{(\ell)}\}_{\ell=1}^L \sim q(\mathbf{z})$
  - Now approximate the expectation as follows

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

- This is basically “weighted” Monte Carlo integration
  - $w^{(\ell)} = \frac{p(\mathbf{z}^{(\ell)})}{q(\mathbf{z}^{(\ell)})}$  denotes the **importance weight** of each sample  $\mathbf{z}^{(\ell)}$

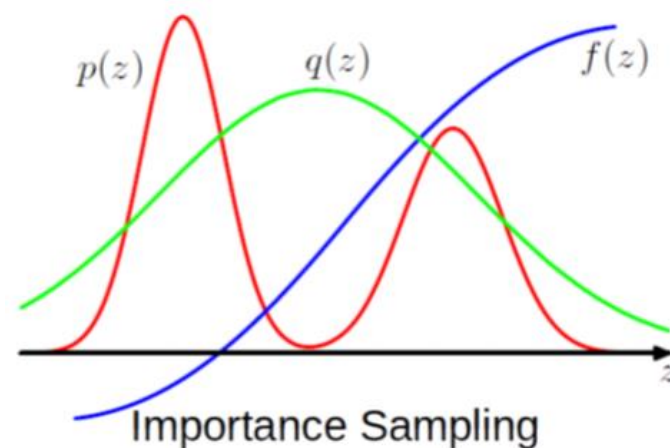
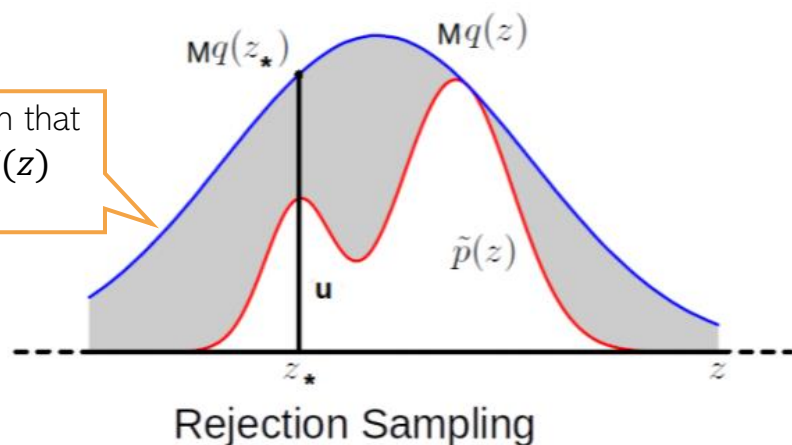
See PRML 11.1.4

- IS works even when we can only evaluate  $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$  up to a prop. constant
- Note: Monte Carlo and Importance Sampling are NOT sampling methods!
  - These are only uses for computing expectations (approximately)



# Limitations of the Basic Methods

- Transformation based methods: Usually limited to drawing from standard distributions
- Rejection Sampling and Importance Sampling: Require good proposal distributions



$$\mathbb{E}[f] \approx \frac{1}{L} \sum_{\ell=1}^L f(z^{(\ell)}) \frac{p(z^{(\ell)})}{q(z^{(\ell)})}$$

Ideally, would like  $q(z)$  to give samples from where  $p(z)$  is large or  $f(z)p(z)$  is large

Difficult to guarantee so if  $z$  is high-dimensional

- In general, difficult to find good prop. distr. especially when  $z$  is high-dim
- More sophisticated sampling methods like MCMC work well in such high-dim spaces



# Markov Chain Monte Carlo (MCMC)

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

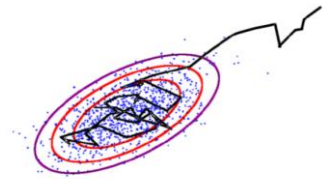
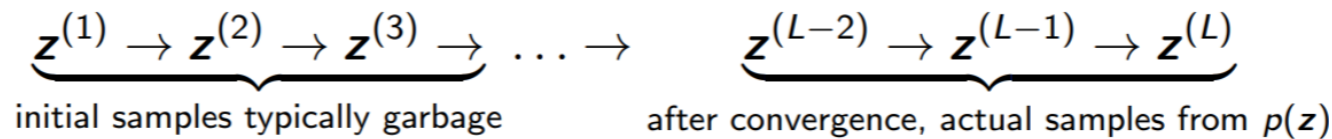
- Goal: Generate samples from some target distribution  $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$

$\mathbf{z}$  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

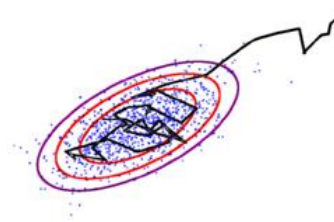
- Use a **proposal distribution**  $q(\mathbf{z}|\mathbf{z}^{(\ell)})$  to generate a candidate sample  $\mathbf{z}_*$
- **Accept/reject**  $\mathbf{z}_*$  as the next sample based on an **acceptance criterion** (will see later)
- If accepted, set  $\mathbf{z}^{(\ell+1)} = \mathbf{z}_*$ . If rejected,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

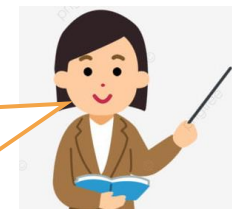


# MCMC: The Basic Scheme



- 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



# 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.,  $\mathcal{N}(\mathbf{z}|\mathbf{z}^{(\tau)}, \sigma^2 \mathbf{I})$
- Suppose the most recently accepted sample in the chain be  $\mathbf{z}^{(\tau)}$
- Draw the next candidate  $\mathbf{z}^*$  from  $q(\mathbf{z}|\mathbf{z}^{(\tau)})$  and compute its acceptance probability

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

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

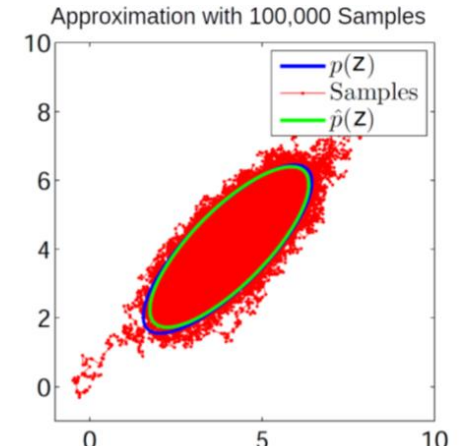
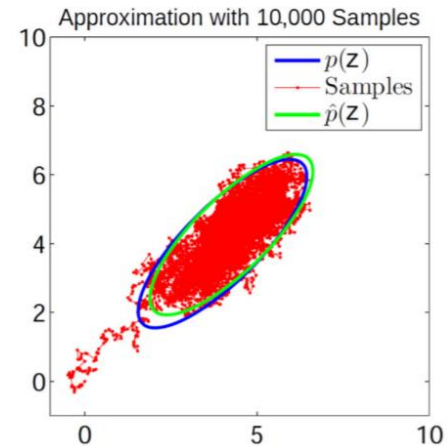
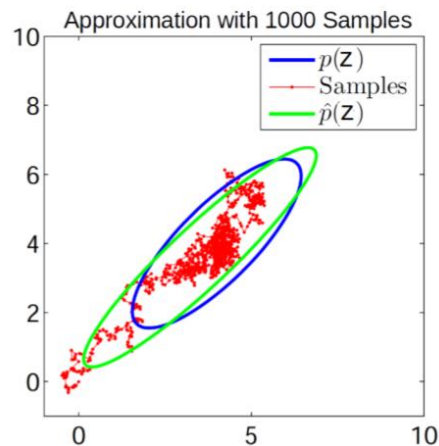
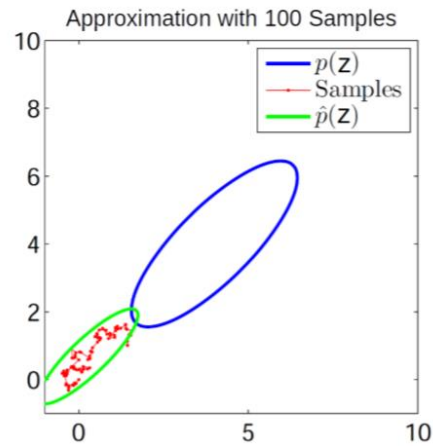
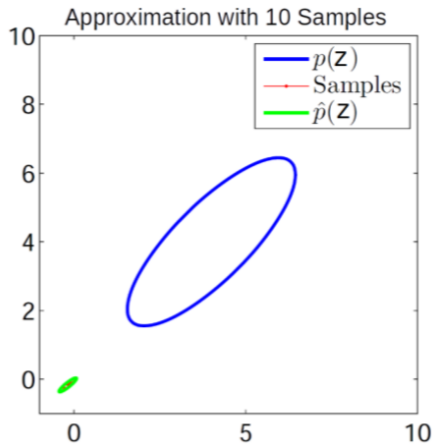
Downweight the probability of acceptance of  $\mathbf{z}^*$  if the proposal itself favors its generation (i.e., if  $q(\mathbf{z}^*|\mathbf{z}^{(\tau)})$  is high), and upweight if it unfavors the generation

- Draw  $u \sim \text{Unif}(0,1)$ . If  $A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) > u$ , then accept  $\mathbf{z}^*$ , set  $\mathbf{z}^\tau = \mathbf{z}^*$ ; otherwise reject  $\mathbf{z}^*$ 
  - This means we are accepting the sample  $\mathbf{z}^*$  with probability  $A(\mathbf{z}^*, \mathbf{z}^{(\tau)})$
- This will generate a chain  $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(T)}$  of samples



# 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  $q(\mathbf{z}^{(t)} | \mathbf{z}^{(t-1)}) = \mathcal{N} \left( \mathbf{z}^{(t-1)}, \begin{bmatrix} 0.01 & 0 \\ 0 & 0.01 \end{bmatrix} \right)$

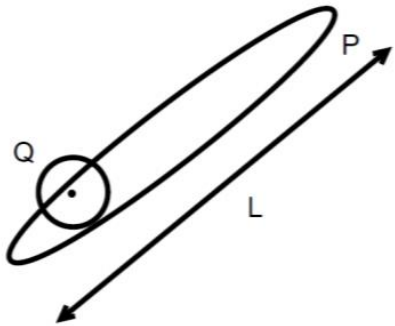


# MH Sampling: Some Comments

- 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)

