# Sampling Methods to Approximate Distributions and Expectations

CS772A: Probabilistic Machine Learning

Piyush Rai

# Approximating a Prob. Distribution using Samples

Can approximate any distribution using a set of randomly drawn samples from it



- The samples can also be used for computing expectations (Monte-Carlo averaging)
- 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.

**CS772A: PML** 

### The Empirical Distribution

- Sampling based approx. can be formally represented using an empirical distribution
- Given L points/samples  $z^{(1)}, z^{(2)}, \dots, z^{(L)}$ , empirical distr. defined by these is





# Sampling: Some Basic Methods

- 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)
- 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 Unif(0, 1), generate  $(z_1, z_2)$  from  $\mathcal{N}(0, \mathbf{I}_2)$ 

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

- Transformation Methods are simple but have limitations
  - Mostly limited to standard distributions and/or distributions with very few variables





p(z) = q(x)

## **Rejection Sampling**

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

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

- Rejection Sampling then works as follows
  - Sample a random variable  $z_*$  from q(z)
  - Sampling a uniform r.v.  $u \sim \text{Unif}[0, Mq(z_*)]$
  - If  $u \leq \widetilde{p}(z_*)$  then accept  $z_*$ , otherwise reject it
- All accepted  $z_*$ 's will be random samples from p(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|accept)

$$p(\operatorname{accept}|z) = \int_{0}^{\tilde{p}(z)} \frac{1}{Mq(z)} du = \frac{\tilde{p}(z)}{Mq(z)}$$

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

$$p(\operatorname{accept}) = \int \frac{\tilde{p}(z)}{M} dz = \frac{Z_{p}}{M}$$

$$p(z|\operatorname{accept}) = \frac{p(z, \operatorname{accept})}{p(\operatorname{accept})} = \frac{\tilde{p}(z)}{Z_{p}} = p(z)$$



### Computing Expectations via Monte Carlo Sampling

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  $p(z): \{z^{(\ell)}\}_{\ell=1}^{L} \sim p(z) \prec$  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 
$$\mathbb{E}[\hat{f}] = \mathbb{E}[f]$$
 and  $\operatorname{var}[\hat{f}] = \frac{1}{L}\operatorname{var}[f] = \frac{1}{L}\mathbb{E}[(f - \mathbb{E}[f])^2]$  Variance in our  
estimate decreases  
as *L* increases

## Computing Expectations via Importance Sampling

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

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

See PRML 11.1.4

CS772A: PML

- This is basically "weighted" Monte Carlo integration
  - $w^{(\ell)} = \frac{p(z^{(\ell)})}{q(z^{(\ell)})}$  denotes the importance weight of each sample  $z^{(\ell)}$
- IS works even when we can only evaluate  $p(z) = \frac{\tilde{p}(z)}{Z_n}$  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



- 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

CS772A: PML



# MCMC: The Basic Scheme

- The chain run infinitely long (i.e., upon convergence) will give ONE sample from  $p({m z})$
- But we usually require several samples to approximate p(z)
- This is done as follows
  - Start the chain at an initial  $m{z}^{(0)}$
  - 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  $z^{(T_1)}$  up to  $T_2$  steps, discard intermediate samples, take last sample  $z^{(T_2)}$ 
    - 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(z)
- Note: Good choices for  $T_1$  and  $T_i T_{i-1}$  (thinning gap) are usually based on heuristics





CS772A: PML

Will treat it as our first sample from p(z)

Requirement for Monte Carlo approximation

MCMC is exact in theory but approximate in practice since

### MCMC: Some Basic Theory

- A first order Markov Chain assumes  $p(\mathbf{z}^{(\ell+1)}|\mathbf{z}^{(1)},...,\mathbf{z}^{(\ell)}) = p(\mathbf{z}^{(\ell+1)}|\mathbf{z}^{(\ell)})$
- 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):  $T_{\ell}(z^{(\ell)} \rightarrow z^{(\ell+1)}) = p(z^{(\ell+1)}|z^{(\ell)})$  the proposal distribution
- TF is a <u>distribution</u> over the values of next state given the value of the current state
- Assuming z is discrete with K possible values, the TF will be  $K \times K$  probability table

Transition probabilities can be defined using a *KxK* table if **z** is a discrete r.v. with *K* possible values



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

