

# Importance Sampling

Sargur N. Srihari  
srihari@cedar.buffalo.edu

# Topics in Monte Carlo Methods

1. Sampling and Monte Carlo Methods
2. Importance Sampling
3. Markov Chain Monte Carlo Methods
4. Gibbs Sampling
5. Mixing between separated modes

# Importance Sampling: Choice of $p(\mathbf{x})$

- Sum of integrand to be computed:

Discrete case

$$s = \sum_{\mathbf{x}} p(\mathbf{x})f(\mathbf{x}) = E_p [f(\mathbf{x})]$$

Continuous case

$$s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x} = E_p [f(\mathbf{x})]$$

- An important step is deciding which part of the integrand has the role of the probability  $p(\mathbf{x})$ 
  - from which we sample  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$
- And which part has role of  $f(\mathbf{x})$  whose expected value (under the probability distribution) is estimated as

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

# Decomposition of Integrand

- In the equation  $s = \sum_x p(\mathbf{x})f(\mathbf{x})$   $s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x}$
- There is no unique decomposition because  $p(\mathbf{x})f(\mathbf{x})$  can be rewritten as

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

– where we now sample from  $q$  and average

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

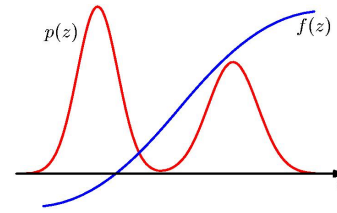
– In many cases the problem  $s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x}$  is specified naturally as expectation of  $f(\mathbf{x})$  given distribution  $p(\mathbf{x})$

- But it may not be optimal in no of samples required

# Difficulty of sampling from $p(\mathbf{x})$

- Principal reason for sampling  $p(\mathbf{x})$  is evaluating expectation of some  $f(\mathbf{x})$

$$E[f] = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$$



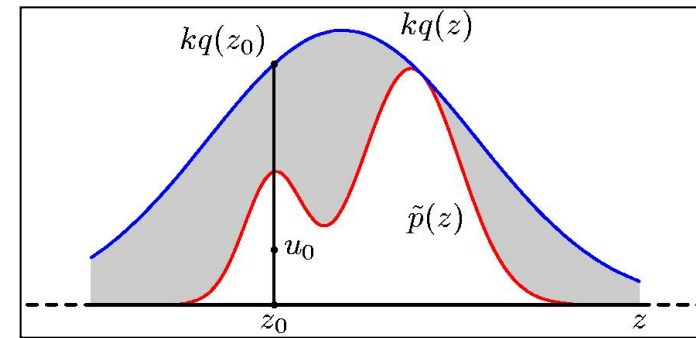
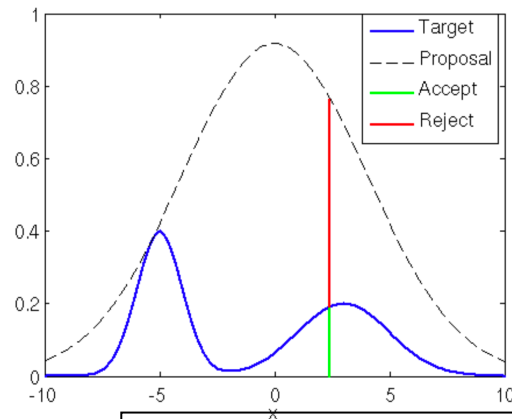
- Given samples  $\mathbf{x}^{(i)}, i=1, \dots, n$ , from  $p(\mathbf{x})$ , the finite sum approximation is

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

- But drawing samples  $p(\mathbf{x})$  may be impractical

# Using a proposal distribution

Rejection sampling uses a proposal distribution



- Determining if sample is in shaded area
- Generate two random numbers
  - $z_0$  from  $q(z)$
  - $u_0$  from uniform distribution  $[0, kq(z_0)]$
- This pair has uniform distribution under the curve of function  $kq(z)$
- If  $u_0 > \tilde{p}(z_0)$  the pair is rejected otherwise it is retained
- Remaining pairs have a uniform distribution under the curve of  $p(z)$  and hence the corresponding  $z$  values are distributed according to  $p(z)$  as desired

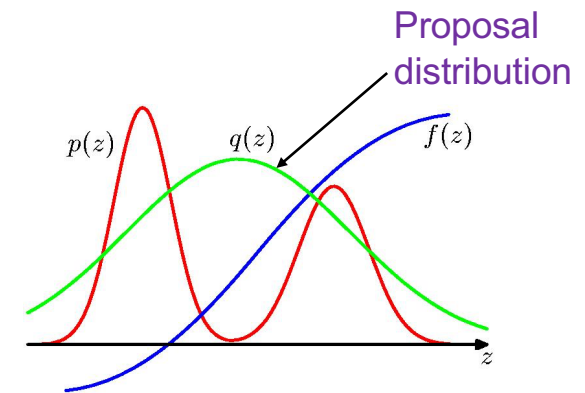
# Importance sampling retains samples

- Importance sampling uses:
  - A proposal distribution– like rejection sampling
    - where samples not matching conditioning are rejected
    - But all samples are retained
  - Assumes that for any  $\mathbf{x}$ ,  $p(\mathbf{x})$  can be evaluated

# Determining Importance weights

- Samples  $\{\mathbf{x}^{(i)}\}$  from simpler distribution  $q(\mathbf{x})$

$$\begin{aligned}
 E[f] &= \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\
 &= \int f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x})d\mathbf{x} \\
 &= \frac{1}{n}\sum_{l=1}^n \frac{p(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})}f(\mathbf{x}^{(i)})
 \end{aligned}$$



- Samples are weighted by ratios

$$r_i = p(\mathbf{x}^{(i)}) / q(\mathbf{x}^{(i)})$$

– Known as importance weights

- Which corrects the bias introduced by wrong distribution

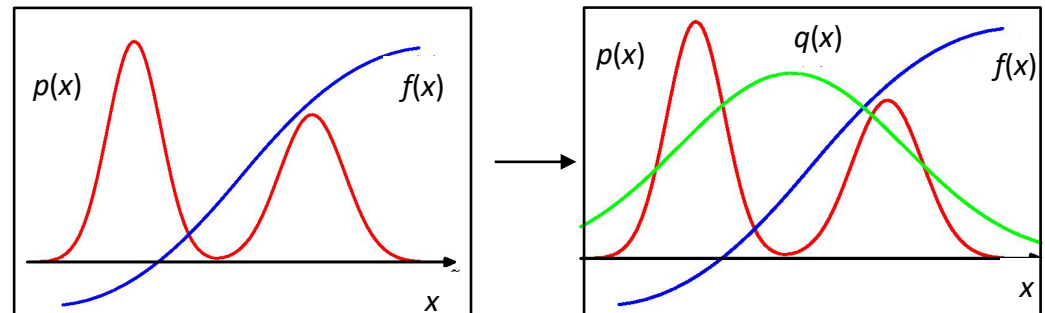


# Importance Sampling

- Better way to compute integrals than rejection
- Sample from a distribution  $q(\mathbf{x})$  and reweigh the samples in a principled way

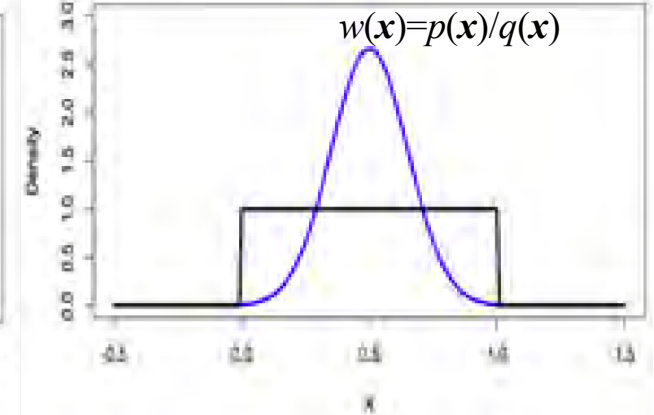
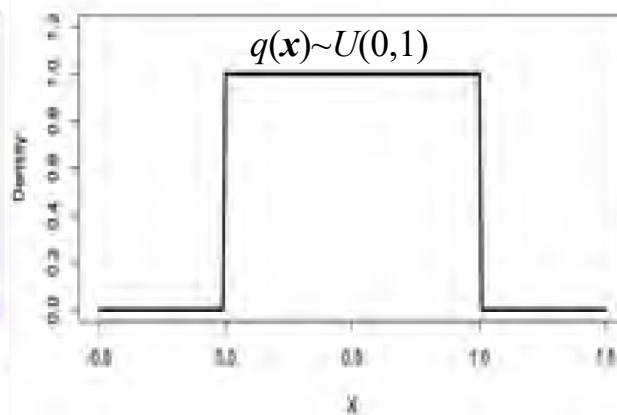
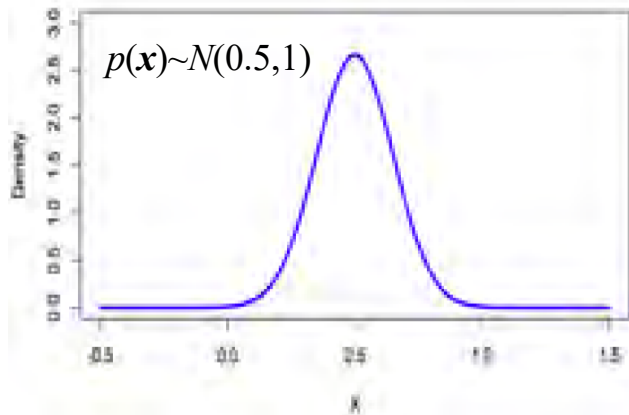
$$\begin{aligned}
 E_{x \sim p}[f(\mathbf{x})] &= \int f(\mathbf{x})p(\mathbf{x}) \\
 &= \int f(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} q(\mathbf{x}) \\
 &= E_{x \sim q}[f(\mathbf{x})w(\mathbf{x})] \\
 &\approx \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^{(m)})w(\mathbf{x}^{(m)})
 \end{aligned}$$

This can be treated as computing expectation of  $p(\mathbf{x})/q(\mathbf{x})$  wrt  $q(\mathbf{x})$



- where  $w(\mathbf{x})=p(\mathbf{x})/q(\mathbf{x})$  and samples  $\mathbf{x}^{(m)}$  drawn from  $q$
- So take samples from  $q(\mathbf{x})$  and reweigh them with  $w(\mathbf{x})$ 
  - But  $p(\mathbf{x})$  is intractable, but there is a solution

# Case of Uniform $q(x)$



If  $\chi = Z \cup E$  then to compute marginal  $p(E=e)$

$$\begin{aligned}
 p(E=e) &= \mathbb{E}_{z \sim p}[p(e|z)] \\
 &= \mathbb{E}_{z \sim q} \left[ p(e|z) \frac{p(z)}{q(z)} \right] \\
 &= \mathbb{E}_{z \sim q} \left[ \frac{p(e, z)}{q(z)} \right] \\
 &= \mathbb{E}_{z \sim q}[w_e(z)] \\
 &\approx \frac{1}{M} \sum_{m=1}^M w_e(z^{(m)})
 \end{aligned}$$

$$p(e) = \sum_z p(e|z)p(z)$$

$$p(e) = \sum_z p(e|z)p(z) \frac{q(z)}{q(z)}$$

where  $w_e(z) = p(e, z)/q(z)$  is the sample weight

$$x = (e, z)$$

Unlike rejection sampling importance sampling will use all samples

# Dealing with intractable $p(\mathbf{x})$

- Importance sampling seems to assume that we are able to compute the importance weight  $w(\mathbf{x}) = p(\mathbf{x}) / q(\mathbf{x})$
- While  $q(\mathbf{x})$  is left as our choice  $p(\mathbf{x})$  may be intractable.
- We can use the unnormalized  $\tilde{p}(\mathbf{x})$  as shown next

# Unnormalized Importance sampling

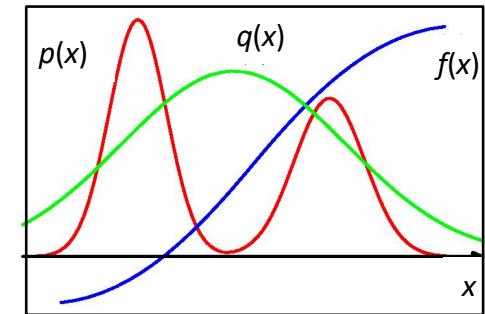
- Define weight using unnormalized distribution

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

- Then

$$E_{x \sim q} [w(\mathbf{x})] = \sum_x q(\mathbf{x}) \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} = \sum_x \tilde{p}(\mathbf{x}) = \alpha$$

- And



$$\begin{aligned} E_{x \sim p} [f(\mathbf{x})] &= \sum_x p(\mathbf{x}) f(\mathbf{x}) \\ &= \sum_x q(\mathbf{x}) f(\mathbf{x}) \frac{p(\mathbf{x})}{q(\mathbf{x})} \\ &= \frac{1}{\alpha} \sum_x q(\mathbf{x}) f(\mathbf{x}) \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} \\ &= \frac{1}{\alpha} E_{x \sim q} [f(\mathbf{x}) w(\mathbf{x})] \\ &= \frac{E_{x \sim q} [f(\mathbf{x}) w(\mathbf{x})]}{E_{x \sim q} [w(\mathbf{x})]} \end{aligned}$$

Therefore, we can determine  $E_{x \sim p} [f(\mathbf{x})]$  using samples drawn from  $q$  as follows:

$$E_{x \sim p} [f(\mathbf{x})] = \frac{\sum_m f(x[m]) w(x[m])}{\sum_m w(x[m])}$$

# Normalized Importance Sampling

- Unnormalized importance sampling unsuitable for conditional probabilities

$$P(X_i = x_i | E = e) = \frac{P(X_i = x_i, E = e)}{P(E = e)}$$

– With unnormalized we could estimate numerator as

$$\begin{aligned} P(X_i = x_i, E = e) &= \sum_z \delta(z) p(e, z) \\ &= \sum_z \delta(z) w_e(z) q(z) \\ &= \mathbb{E}_{z \sim q} [\delta(z) w_e(z)] \\ &\approx \frac{1}{M} \sum_{m=1}^M \delta(z^{(m)}) w_e(z^{(m)}) \end{aligned}$$

where  $\delta(z) = \begin{cases} 1 & \text{if } z \text{ is consistent with } X_i = x_i \\ 0 & \text{otherwise} \end{cases}$   
and  $w_e(z) = p(e, z) / q(z)$

– and denominator from samples as

$$P(E = e) \approx \frac{1}{M} \sum_{m=1}^M w_e(z^{(m)})$$

Which is shown on Slide 16

- Separate Numerator & Denominator compounds errors
  - » If Numerator is an underestimate, Denominator is an overestimate, overall result is a severe underestimate

# Using same set of samples

- If we use the same set of  $M$  samples  $z^1, \dots, z^M \sim q$ 
  - For both numerator and denominator, we avoid this issue of compounding errors
  - Final form of normalized importance sampling is

$$\hat{P}(X_i = x_i | E = e) = \frac{\frac{1}{T} \sum_{t=1}^T \delta(z^t) w_e(z^t)}{\frac{1}{T} \sum_{t=1}^T w_e(z^t)} \quad \text{where} \quad \begin{cases} \delta(z) = \begin{cases} 1 & \text{if } z \text{ is consistent with } X_i = x_i \\ 0 & \text{otherwise} \end{cases} \\ \text{and } w_e(z) = p(e, z) / q(z) \end{cases}$$

- However, it is biased, for  $M=1$

$$\mathbb{E}_{z \sim q}[\hat{P}(X_i = x_i | E = e)] = \mathbb{E}_{z \sim q}[\delta(z)] \neq P(X_i = x_i | E = e)$$

- Fortunately both numerator and denominator are unbiased so the result is asymptotically unbiased

$$\lim_{T \rightarrow \infty} \mathbb{E}_{z \sim q}[\hat{P}(X_i = x_i | E = e)] = P(X_i = x_i | E = e).$$

# Derivation of Optimal $q^*(\mathbf{x})$

- Any Monte Carlo estimator

$$\hat{s}_p = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim p}^n f(\mathbf{x}^{(i)})$$

- can be transformed into an importance sampling estimator

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

using

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

- It can be readily seen that the expected value of the estimator does not depend on  $q$ :  $E_q[\hat{s}_q] = E_q[\hat{s}_p] = s$

- The variance of an importance sampling estimator is sensitive to the choice of  $q$ . The variance is

- The minimum variance occurs when  $q$  is

$$q^*(\mathbf{x}) = \frac{p(\mathbf{x}) |f(\mathbf{x})|}{Z}$$

$$Var[\hat{s}_p] = Var\left[\frac{1}{n} \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right]$$

- where  $Z$  is the normalization constant chosen so that  $q^*(\mathbf{x})$  sums or integrates to one

# Choice of suboptimal $q(\mathbf{x})$

- Any choice of  $q$  is valid and  $q^*$  is the optimal one (yields minimum variance)
- Sampling from  $q^*$  is usually infeasible
- Other choices of  $q$  can be feasible, reducing variance somewhat



# Biased Importance Sampling (BIS)

- Another approach is BIS
  - Has advantage of not requiring normalized  $p$  or  $q$
  - For discrete variables, BIS estimator is

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

where  $\hat{p}$  and  $\hat{q}$  are unnormalized forms of  $p$  and  $q$  and  $\mathbf{x}^{(i)}$  are samples from  $q$

# Effect of choice of $q$

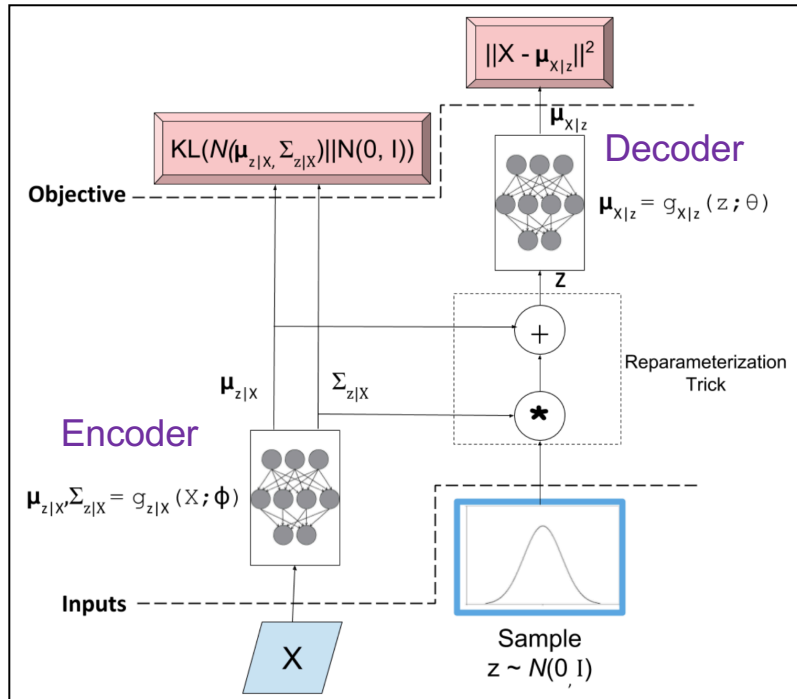
- Good  $q \rightarrow$  efficient Monte Carlo estimation
- Poor choice of  $q \rightarrow$  efficiency much worse
  - Looking at  $\text{Var}[\hat{s}_p] = \text{Var}\left[\frac{1}{n} \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right]$  if there are samples for which  $\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}$  is large then the variance of the estimator can get large
    - Happens when  $q(\mathbf{x})$  is tiny while neither  $p(\mathbf{x})$  or  $f(\mathbf{x})$  is small enough to cancel it
  - The  $q$  distribution is usually chosen to be a very simple distribution so that it is easy to sample from
    - When  $\mathbf{x}$  is high dimensional this simplicity causes it to match  $p$  or  $|f|$  poorly
    - Very small or very large ratios are possible when  $\mathbf{x}$  is high dimensional

# Importance sampling in Deep learning

- Used in many ML algorithms incl. deep learning
- Examples
  - To accelerate training in neural language models with large vocabulary
    - Other neural nets with large no of outputs
  - Estimate partition function (normalize prob. distribution)
  - Estimate log-likelihood in deep directed models such as variational autoencoder
  - Estimate gradient in SGD where most of the cost comes from a small no of misclassified samples
    - Sampling more difficult examples can reduce variance of gradient

# Importance Sampling and VAE

To evaluate the goodness of a VAE we decide to use Bayesian Model comparison



$$\Pr(M|D) = \frac{\Pr(D|M) \Pr(M)}{\Pr(D)}$$

Thus for each model  $M$  we need

$$P_M(X) = \int p_M(X|z)p(z)dz \approx \frac{1}{N} \sum_{i=1}^N p_M(x|z_i), \quad z_i \sim \mathcal{N}(0, 1)$$

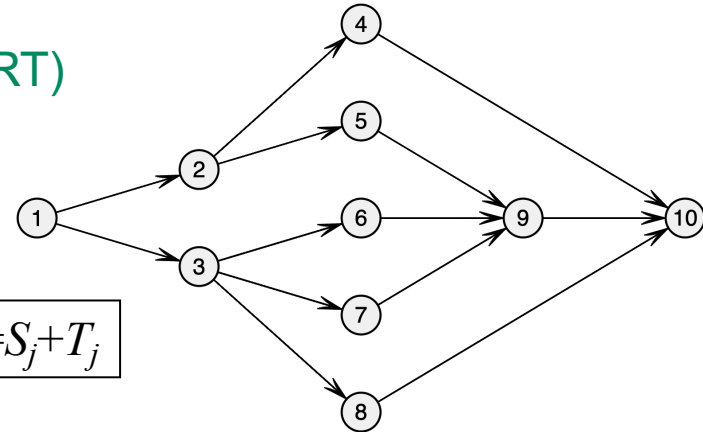
Use importance sampling with  $q$  being scaled and shifted Gaussians from our encoder

$$P_M(X) = \int p_M(X|z)p(z)dz \approx \frac{1}{N} \sum_{i=1}^N \frac{p_M(X|z_i)p_{\mathcal{N}(0,1)}(z_i)}{q_{\mathcal{N}(\mu(X), \sigma(X))}(z_i)}, \quad z_i \sim \mathcal{N}(\mu(X), \sigma(X))$$

# Importance Sampling in PERT

## Project Evaluation and Review Technique (PERT)

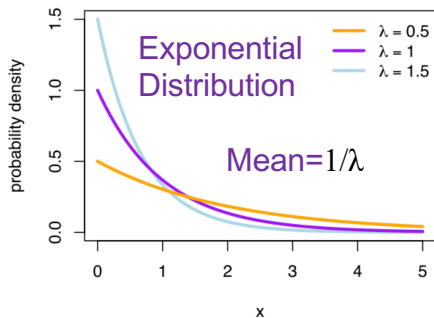
$j$	Task	Predecessors	Duration $\bar{T}_j$
1	Planning	None	4
2	Database Design	1	4
3	Module Layout	1	2
4	Database Capture	2	5
5	Database Interface	2	2
6	Input Module	3	3
7	Output Module	3	2
8	GUI Structure	3	3
9	I/O Interface Implementation	5,6,7	2
10	Final Testing	4,8,9	2



$$E_j = S_j + T_j$$

If every task takes time in table, then  $E_{10}=15$

If  $T_j$  are independent exponentially distributed with means as given, then completion time is a random variable with mean  $E_{10}=18$  and a long tail to the right



If there is a severe penalty for  $E_{10} > 70$ , direct simulation with  $n=10,000$  indicates  $E_{10} > 70$  happened only 2 times

Importance sampling gets a better estimate of this probability

$$\hat{\mu}_\lambda = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{E_{i,10} > 70\} \prod_{j=1}^{10} \frac{\exp(-T_{ij}/\theta_j)/\theta_j}{\exp(-T_{ij}/\lambda_j)/\lambda_j}$$

with  $\lambda = 4\theta$  and  $n = 10,000$  gave  $\hat{\mu}_\lambda = 2.02 \times 10^{-5}$