

# **Probabilistic Graphical Models**

**Monte Carlo Sampling**

# Sampling-Based (Monte Carlo) Estimation

Example: Estimate  $p$  for a Bernoulli distribution  $P$

Toss a coin  $m$  times, compute  $\hat{p} = E[X] = \frac{1}{M} \sum X_i$

Theoretical expectation:  $p = E[X] = \sum_x xP(x)$

Empirical Expectation:  $\hat{p} = \hat{E}[X] \approx \frac{1}{M} \sum X_i$

Estimator is unbiased

$$E[\hat{p}] = p$$

Asymptotics:  $\hat{p} \rightarrow p, \hat{p} \sim \text{Norm}(p, \frac{\sigma^2}{n})$

Daphne Koller

# Sampling-Based (Monte Carlo) Estimation

Idea: Distribution  $P$ , you want to compute  $E[f(\mathbf{X})]$

Exact/Approximate Inference is very hard

Sample  $D = \{\mathbf{x}[1], \dots, \mathbf{x}[M]\}$  (i.i.d)

For a distribution  $P$ , function  $f$ :

Estimate  $\widehat{E}_P[f(\mathbf{X})] = \frac{1}{M} \sum f(\mathbf{x}[m])$

# Properties

Estimator is unbiased:  $E[\hat{f}] = E[\widehat{E}_D[f(X)]] = E_P[f(X)]$

Variance of the estimator:  $\sigma_{\hat{f}}^2 = \sigma_f^2/n$

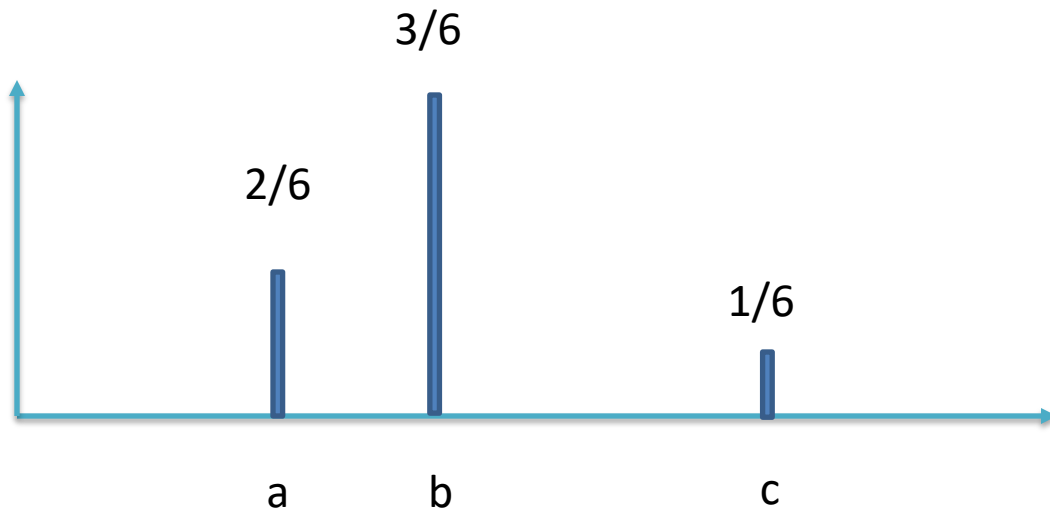
Weak Law of Large Numbers:

$$\lim_{n \rightarrow \infty} P(|\hat{f} - E_P[f(X)]| < \epsilon) = 1$$

CLT:  $\hat{f} \sim N(E_P[f(X)], \sigma_{\hat{f}}^2)$

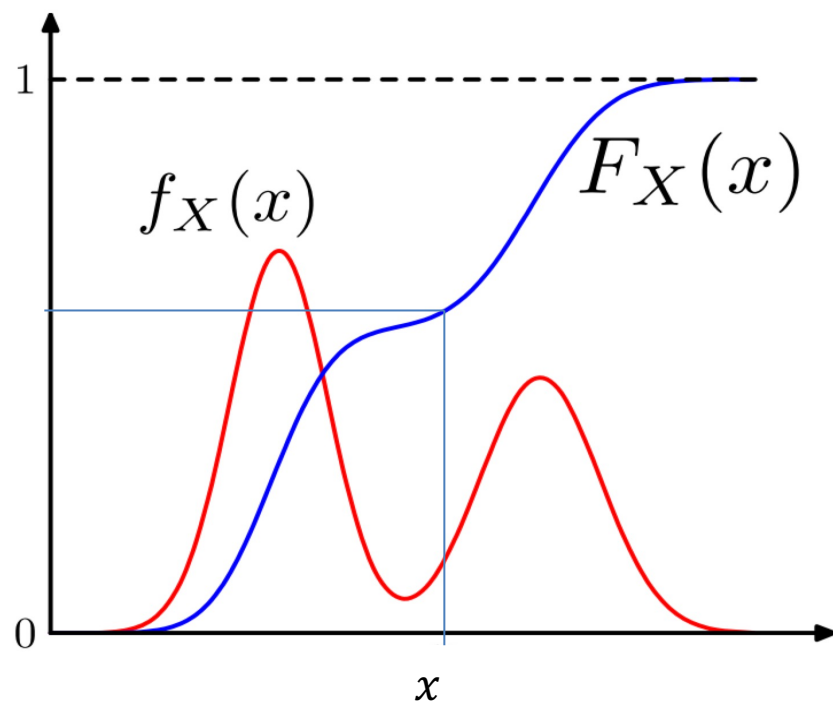
# How do we sample?

- Computers can produce random numbers from a uniform distribution.
- Inverse CDF transform-based sampling



Let  $U_i = F(x_i)$ . Then  $U_i \sim \text{Uniform}(0, 1)$

# How do we sample?



Let  $F$  be a strictly monotonic cumulative distribution function, and let  $F^{-1}$  be its inverse function

Claim: If  $U$  is a uniform random variable on  $[0,1]$  then  $F^{-1}(U)$  has  $F$  as its CDF.

Proof:

$$\begin{aligned}\Pr(F^{-1}(U) \leq x) \\ &= \Pr(U \leq F(x)) \\ &= F(x) \text{ (because } \Pr(U \leq u) = u, \text{ when } U \text{ is uniform on } [0,1])\end{aligned}$$

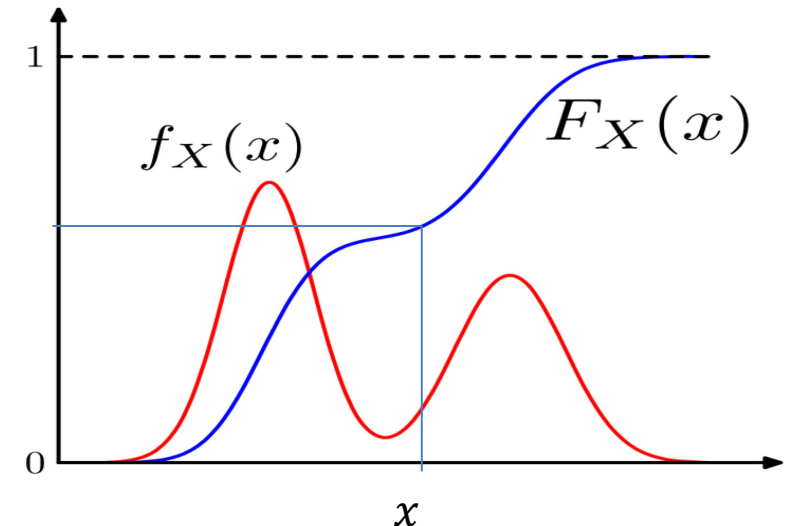
Let  $U_i = F(x_i)$ . Then  $U_i \sim \text{Uniform}(0, 1)$

$$\begin{aligned}P(X_i \leq x) \\ &= P(F_x^{-1}(u_i) \leq x) = F(x)\end{aligned}$$

# How do we sample?

1. Generate a random number  $u$  from the standard uniform distribution in the interval  $[0,1]$ , i.e. from  $U \sim \text{Unif}[0,1]$ .
2. Find the generalized inverse of the desired CDF, i.e.  $F_X^{-1}(u)$ .
3. Compute  $X'(u) = F_X^{-1}(u)$ . The computed random variable  $X'(U)$  has distribution  $F_X$  and thereby the same law as  $X$ .

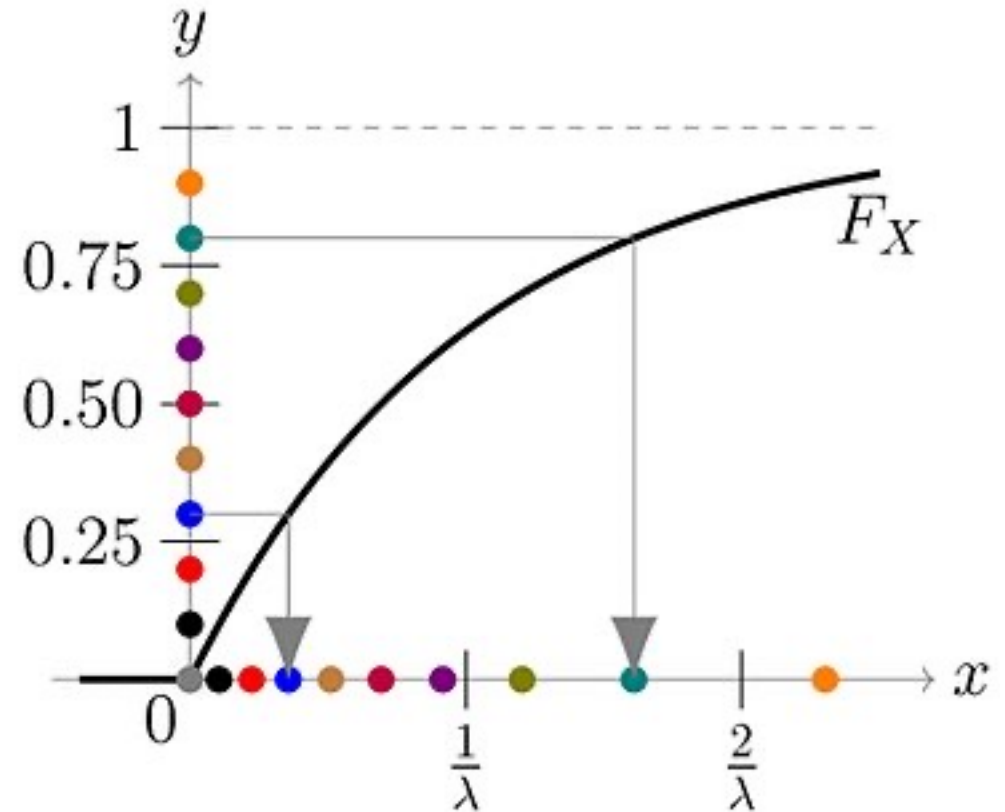
Let  $U_i = F(x_i)$ . Then  $U_i \sim \text{Uniform}(0, 1)$



# Example: Exponential Distribution

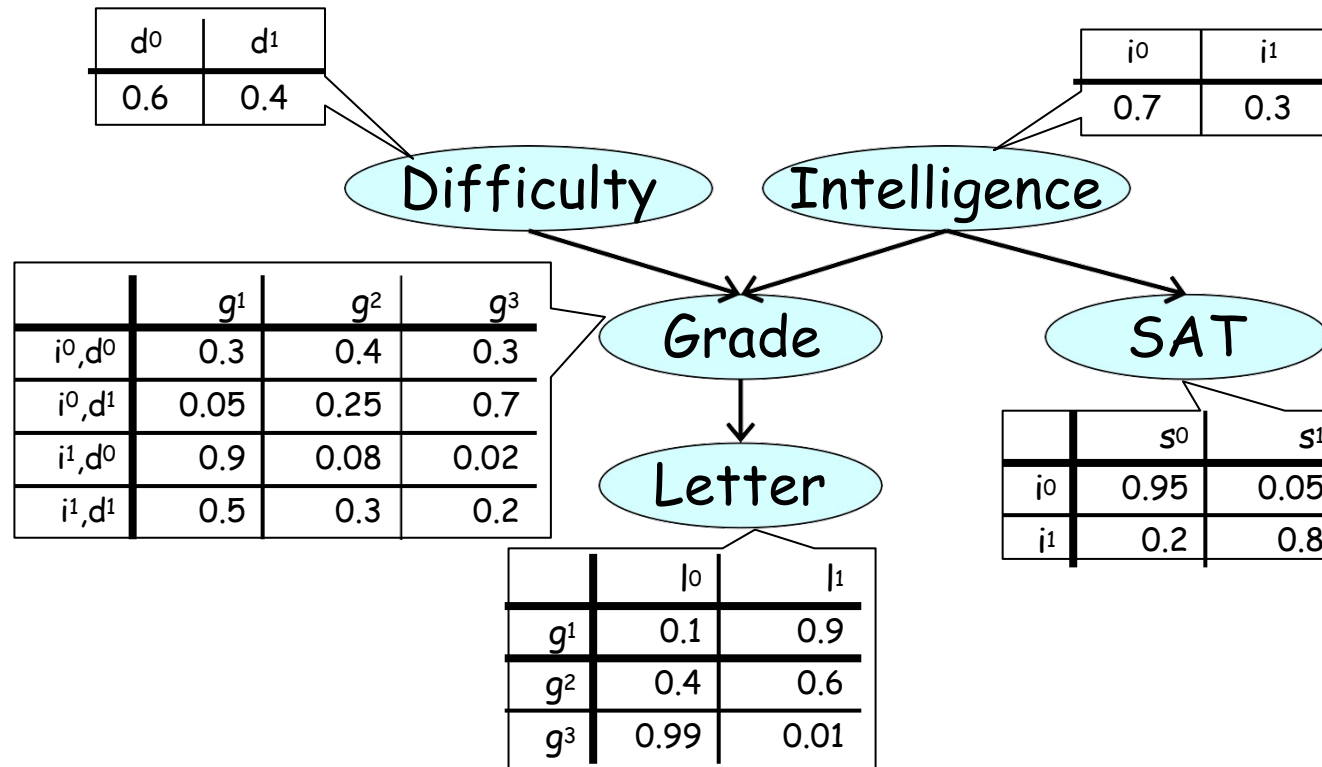
- Draw some  $u_i$  from a  $U \sim \text{Unif}(0,1)$
- Compute  $x_i = F_X^{-1}(u_i) = -\frac{1}{\lambda} \ln(1 - u_i)$
- This  $x_i$  has exponential distribution.

How about  $F(x) = 1 - \exp(-\sqrt{x})$





# Forward sampling from a BN



Daphne Koller

# Sampling-Based (Monte Carlo) Estimation

Idea: Distribution  $P$ , you want to compute  $P(\mathbf{X} = \mathbf{x}) = E[I(\mathbf{X} = \mathbf{x})]$

Exact/Approximate Inference is very hard

Sample  $D = \{\mathbf{x}[1], \dots, \mathbf{x}[M]\}$  (i.i.d)

Estimate  $\hat{P}(\mathbf{X} = \mathbf{x}) = \frac{1}{M} \sum I(\mathbf{X} = \mathbf{x})$

# Sampling-Based (Monte Carlo) Estimation

Hoeffding Bound:

$$P_{\mathcal{D}}(T_{\mathcal{D}} \notin [p - \epsilon, p + \epsilon]) \leq 2e^{-2M\epsilon^2}$$

$$T_{\mathcal{D}} = \frac{1}{M} \sum_{m=1}^M X[m]$$

For additive bound  $\epsilon$  on error with probability  $> 1-\delta$ :

$$M \geq \frac{\ln(2/\delta)}{2\epsilon^2}$$

Chernoff Bound:

$$P_{\mathcal{D}}(T_{\mathcal{D}} \notin [p(1 - \epsilon), p(1 + \epsilon)]) \leq 2e^{-Mp\epsilon^2/3}$$

For multiplicative bound  $\epsilon$  on error with probability  $> 1-\delta$ :

$$M \geq 3 \frac{\ln(2/\delta)}{p\epsilon^2}$$

Daphne Koller

# Forward Sampling for Querying

- Goal: Estimate  $P(Y=y)$ 
  - Generate samples from BN
  - Compute fraction where  $Y=y$

For additive bound  $\epsilon$  on error with probability  $> 1-\delta$

$$M \geq \frac{\ln(2/\delta)}{2\epsilon^2}$$

For multiplicative bound  $\epsilon$  on error with probability  $> 1-\delta$

$$M \geq 3 \frac{\ln(2/\delta)}{P(y)\epsilon^2}$$

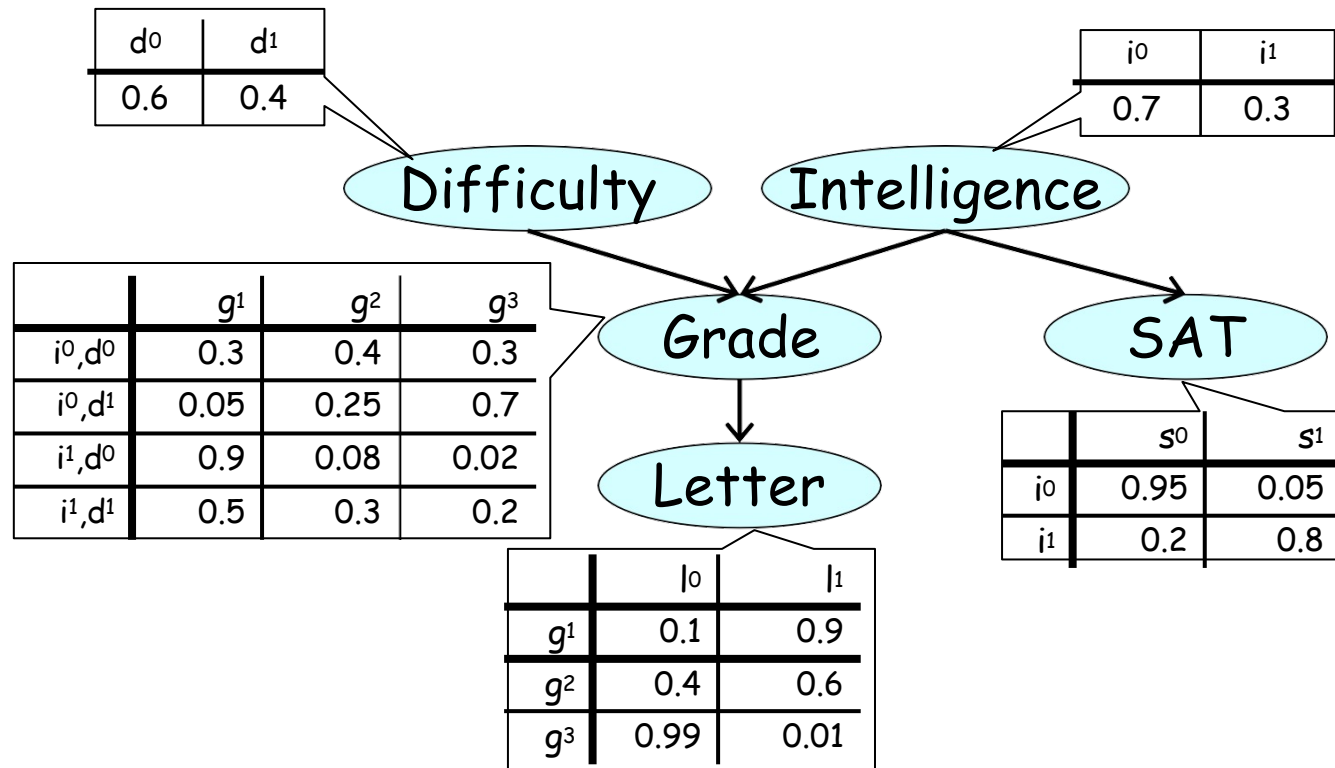
# Queries with Evidence

- Goal: Estimate  $P(Y=y \mid \mathbf{E}=\mathbf{e})$
- Rejection sampling algorithm
  - Generate samples from BN
  - Throw away all those where  $\mathbf{E} \neq \mathbf{e}$
  - Compute fraction where  $Y=y$

Expected fraction of samples kept  $\sim P(\mathbf{e})$

# samples needed rows exponentially with # of observed variables

# Likelihood Weighting Sampling



Example: Force  $s^1$

$$P(s^1 | i^0) = 0.05$$

$$P(s^1 | i^1) = 0.8$$

Example: Force  $s^1, l^0$ ,  
sample  $d^1, i^0, g^2$

$$P(l^0 | g^2) = 0.4$$

$$P(s^1 | i^0) = 0.05$$

# Likelihood Weighting Sampling

---

**Algorithm 12.2 Likelihood-weighted particle generation**

---

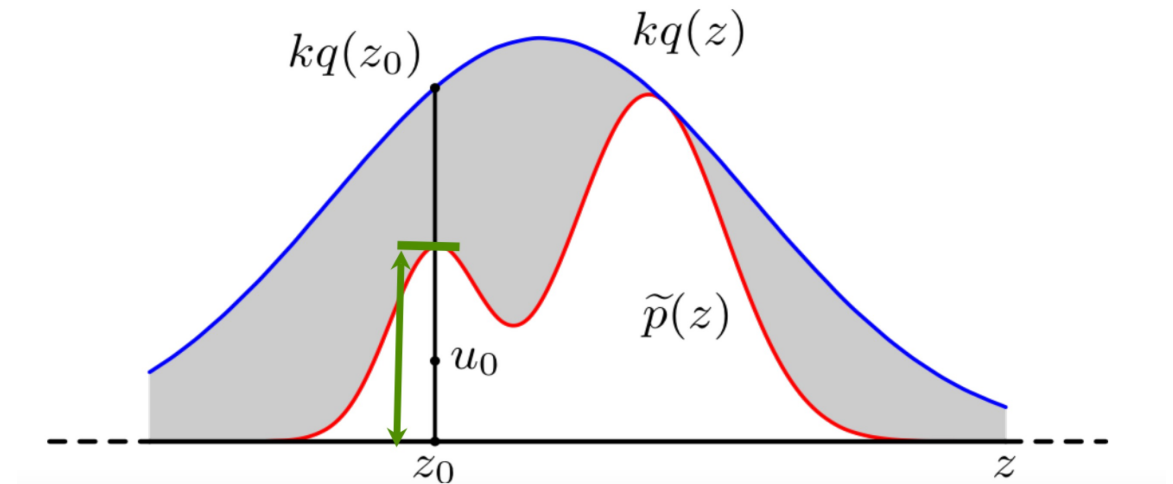
```
Procedure LW-Sample (  
     $\mathcal{B}$ , // Bayesian network over  $\mathcal{X}$   
     $Z = z$  // Event in the network  
)  
1  Let  $X_1, \dots, X_n$  be a topological ordering of  $\mathcal{X}$   
2   $w \leftarrow 1$   
3  for  $i = 1, \dots, n$   
4       $u_i \leftarrow x \langle \text{Pa}_{X_i} \rangle$  // Assignment to  $\text{Pa}_{X_i}$  in  $x_1, \dots, x_{i-1}$   
5      if  $X_i \notin Z$  then  
6          Sample  $x_i$  from  $P(X_i \mid u_i)$   
7      else  
8           $x_i \leftarrow z \langle X_i \rangle$  // Assignment to  $X_i$  in  $z$   
9           $w \leftarrow w \cdot P(x_i \mid u_i)$  // Multiply weight by probability of desired value  
10 return  $(x_1, \dots, x_n), w$ 
```

---

$$\widehat{P}_D(\mathbf{y} \mid \mathbf{e}) = \frac{\sum_{m=1}^M w[m] 1\{\mathbf{y}[m] = \mathbf{y}\}}{\sum_{m=1}^M w[m]}$$

# General Framework: Rejection Sampling

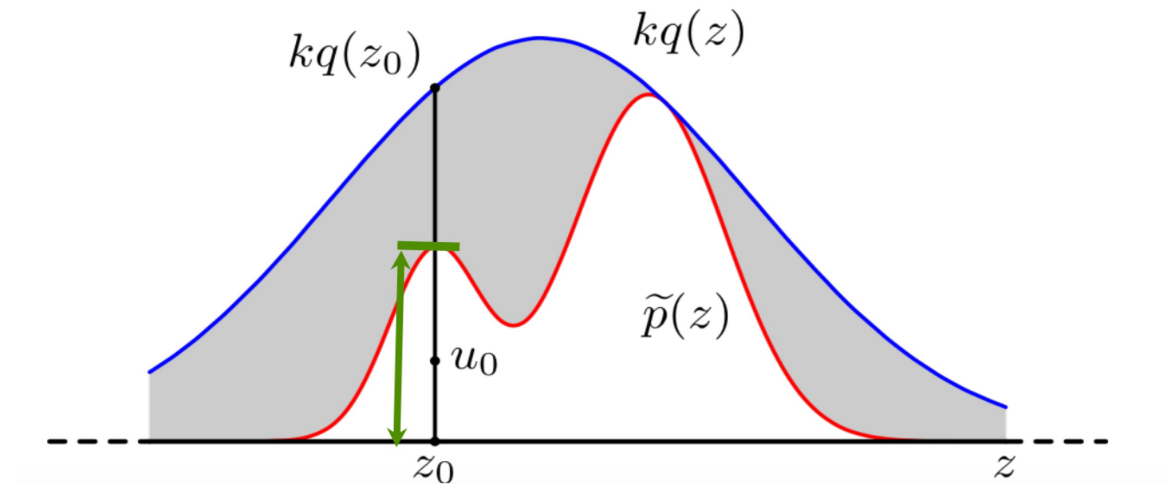
- Access to easy-to-sample distribution  $Q(x)$
- Constant  $k$  such that  $P(x) \leq kQ(x)$
- Algorithm
  1. Sample from  $Q(x)$
  2. Keep samples in proportion to  $\frac{P(x)}{k \cdot Q(z)}$  and reject the rest





# General Framework: Rejection Sampling

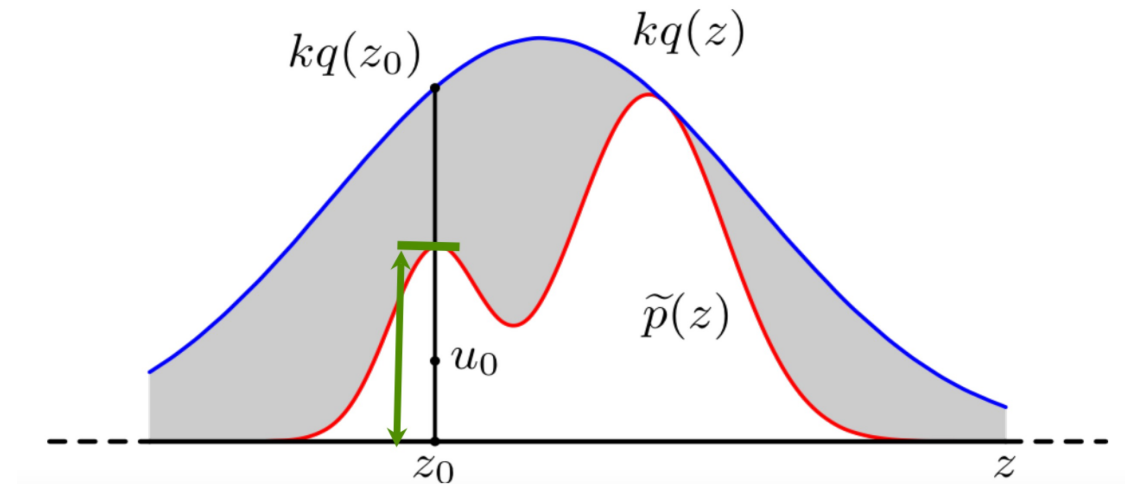
- Access to easy-to-sample distribution  $Q(x)$
- Constant  $k$  such that  $P(x) \leq kQ(x)$
- Algorithm
  1. Sample  $x_i \sim Q(x)$
  2. Sample  $y_i \sim \text{Uni}[0, kQ(x_i)]$
  3. If  $P(x_i) \leq y_i$  keep, otherwise reject.



# General Framework: Rejection Sampling

- Access to easy-to-sample distribution  $Q(x)$
- Constant  $k$  such that  $\tilde{P}(x) \leq kQ(x)$
- Algorithm
  1. Sample  $x_i \sim Q(x)$
  2. Sample  $y_i \sim \text{Uni}[0, kQ(x_i)]$
  3. If  $\tilde{P}(x_i) \leq y_i$  keep, otherwise reject.

Example Uses Gaussian proposal  $q$  to draw samples from multimodal distribution  $p$

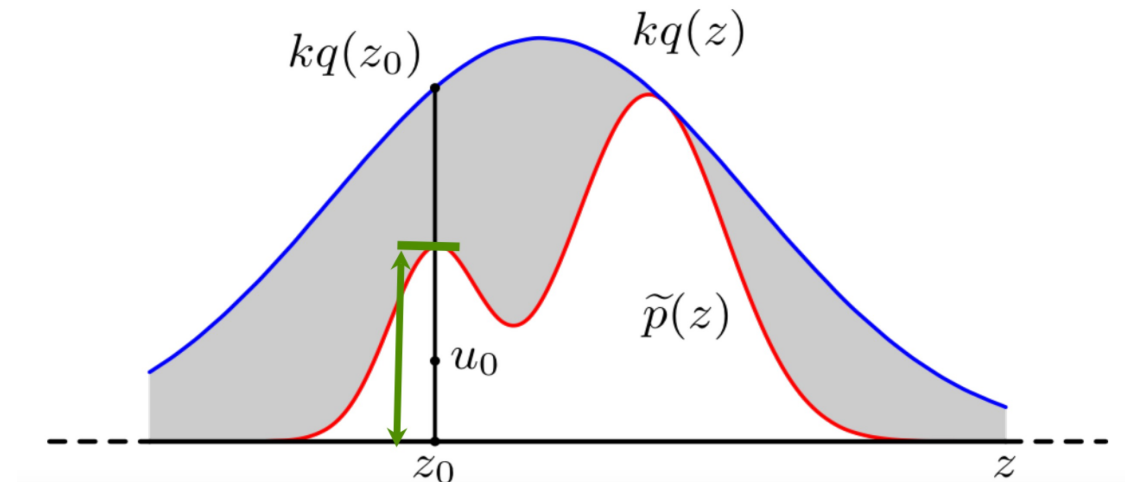


# General Framework: Rejection Sampling

- Access to easy-to-sample distribution  $Q(x)$
- Constant  $k$  such that  $\tilde{P}(x) \leq kQ(x)$
- Algorithm
  1. Sample  $x_i \sim Q(x)$
  2. Sample  $y_i \sim \text{Uni}[0, kQ(x_i)]$
  3. If  $\tilde{P}(x_i) \leq y_i$  keep, otherwise reject.

Can be very wasteful  
Needs to be bounded

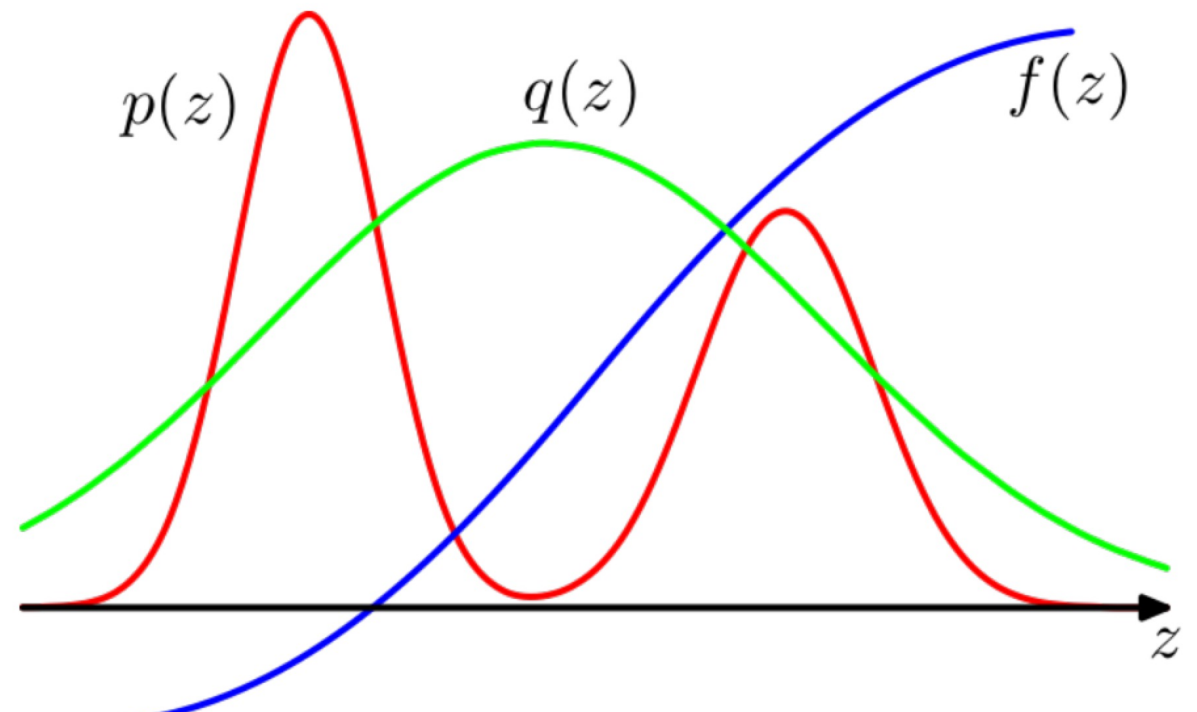
Example Uses Gaussian proposal  $q$  to draw samples from multimodal distribution  $p$



# Importance Sampling (unnormalized)

- Access to easy-to-sample distribution  $Q(x)$
- $Q(x) > 0$  whenever  $P(X) > 0$
- You want to compute  $E[f(X)]$

- Idea: Sample from  $Q(x)$
- Weigh the sample proportionally to
$$w(x) = \frac{P(x)}{Q(x)}$$

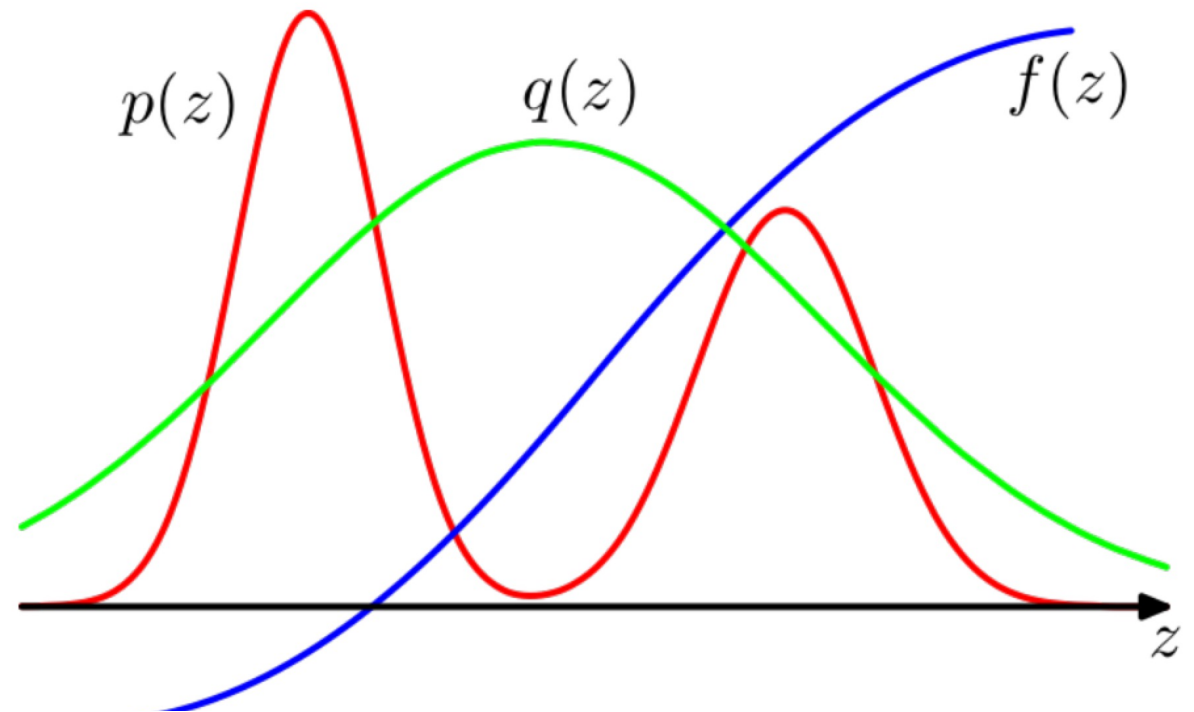


# Importance Sampling (unnormalized)

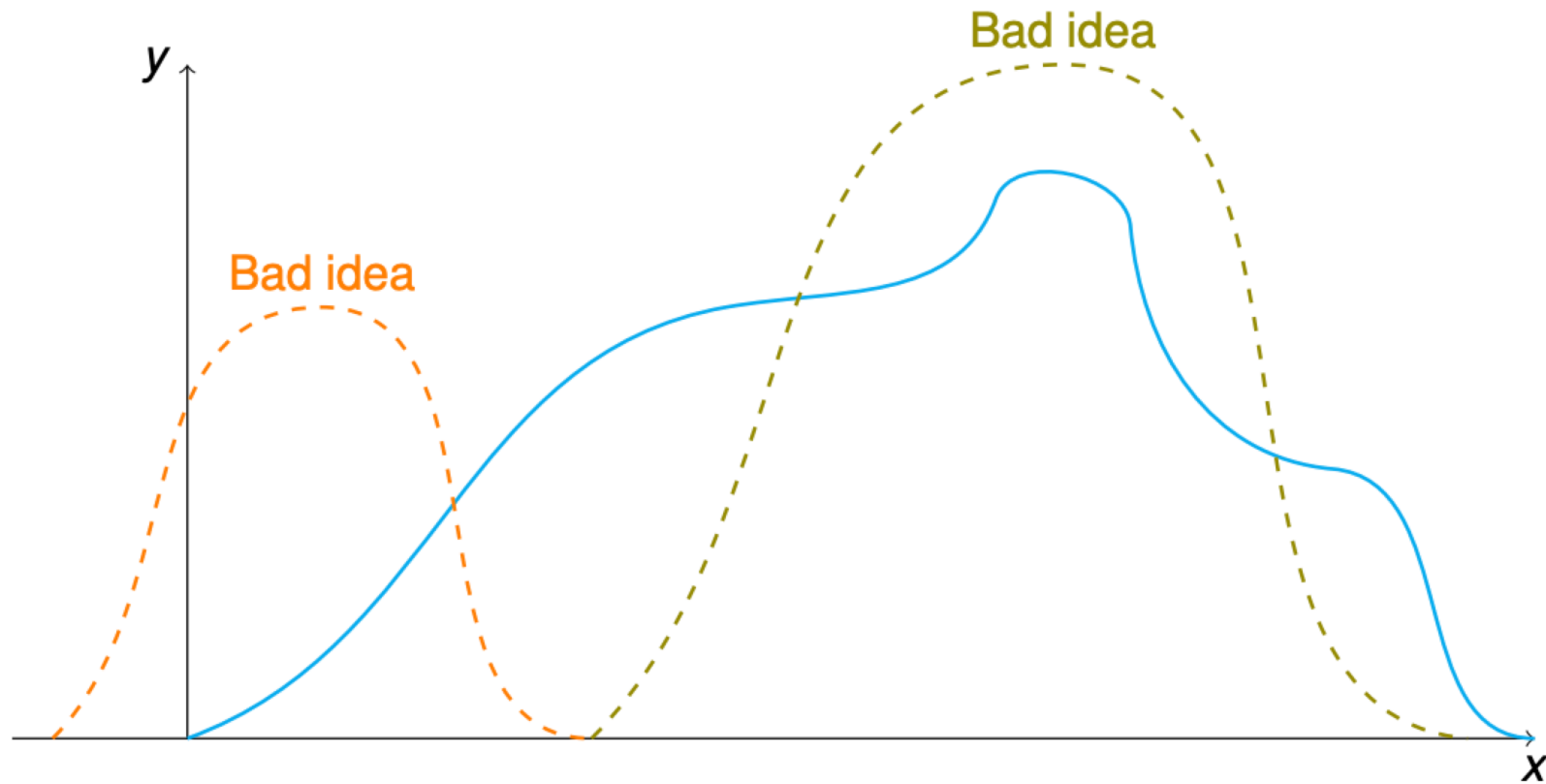
$$E_{P(X)}[f(X)] = E_{Q(X)} \left[ f(X) \frac{P(X)}{Q(X)} \right]$$

$$\approx \frac{1}{M} \sum_{m=1}^M \frac{p(x^{(m)})}{q(x^{(m)})} f(x^{(m)})$$

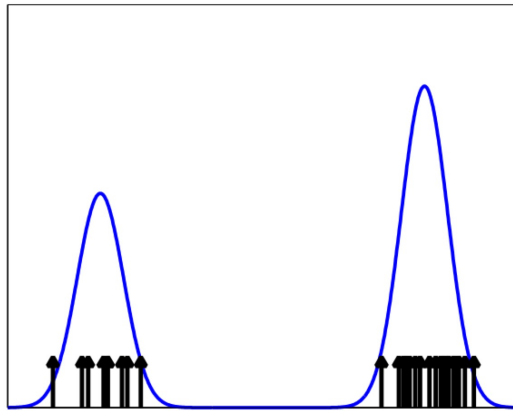
Where  $x_m$  are sampled from  $q(x)$



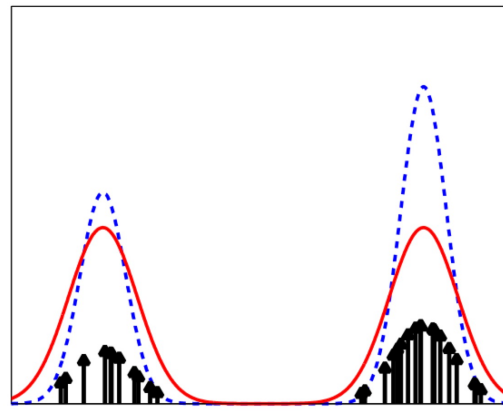
# Selecting Good Proposals



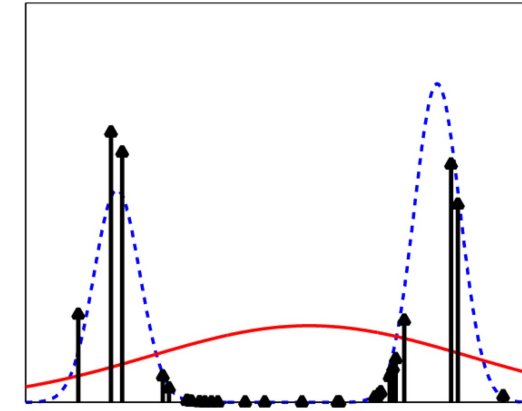
# Selecting Good Proposals



*Target Distribution*



*Good Proposal*

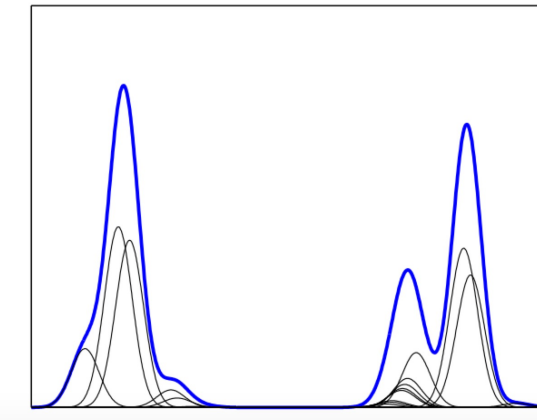
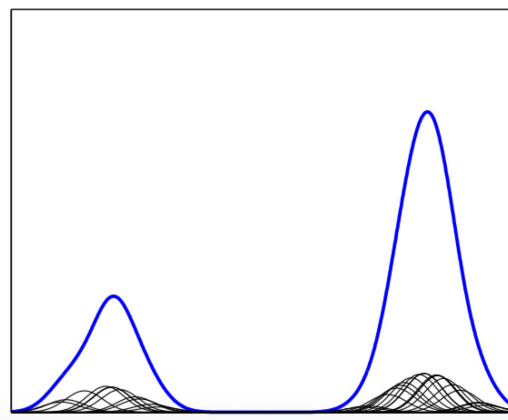
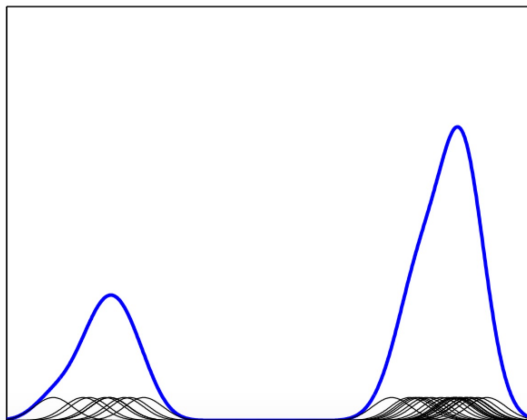


*Poor Proposal*

*Kernel or Parzen window estimators  
interpolate to predict density:*

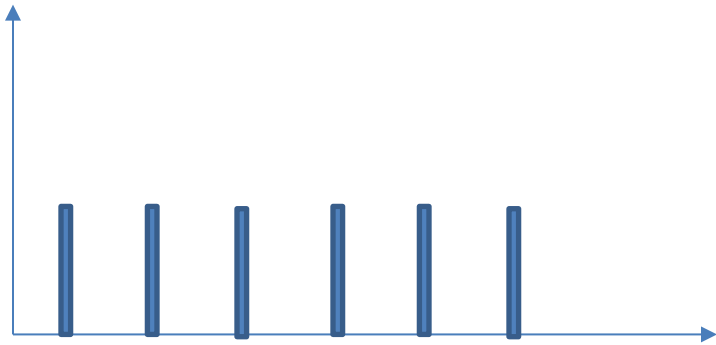
$$\hat{p}(x) = \sum_{\ell=1}^L w^{(\ell)} \mathcal{N}(x; x^{(\ell)}, \Lambda)$$

$$w^{(\ell)} \propto \frac{p(x^{(\ell)})}{q(x^{(\ell)})}$$



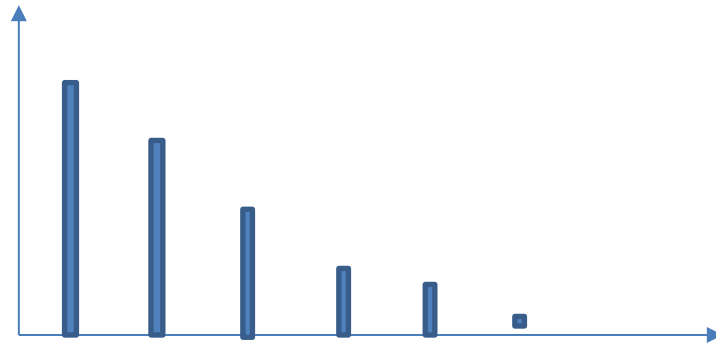
# Example:

Sample from a fair die to simulate data from an unfair die



$$E[X] = 3.5$$

$$P(X = 1) = 1/6$$



$$E[X] = 1.9$$

$$P(X = 1) = 4/6$$



# Bias/Variance of the importance sampling estimator

For data sets  $\mathcal{D}$  sampled from  $Q$ , we have that:

$$\mathbb{E}_{\mathcal{D}}[\hat{\mathbb{E}}_{\mathcal{D}}(f)] = \mathbf{E}_{Q(\mathbf{X})}[f(\mathbf{X})w(\mathbf{X})] = \mathbb{E}_{P(\mathbf{X})}[f(\mathbf{X})]$$

We can also estimate the distribution of this estimator around its mean. Letting  $\epsilon_{\mathcal{D}} = \hat{\mathbb{E}}_{\mathcal{D}}(f) - \mathbb{E}_P[f(x)]$ , we have that, since  $M \rightarrow \infty$ :

$$\mathbb{E}_{\mathcal{D}}[\epsilon_{\mathcal{D}}] \sim \mathcal{N}(0; \sigma_Q^2/M)$$

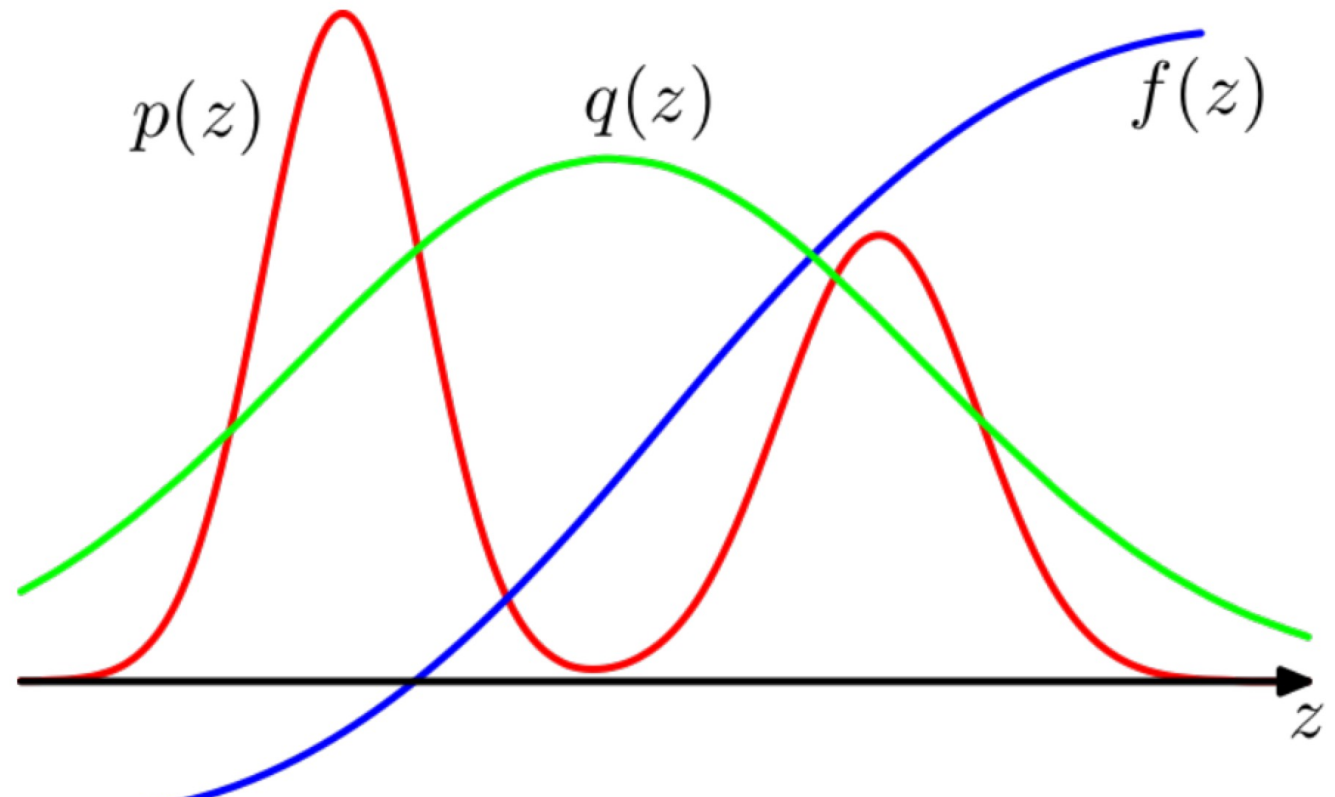
where

$$\begin{aligned}\sigma_Q^2 &= \mathbf{E}_{Q(\mathbf{X})}[(f(\mathbf{X})w(\mathbf{X}))^2] - \mathbb{E}_{Q(\mathbf{X})}[f(\mathbf{X})w(\mathbf{X})]^2 \\ &= \mathbf{E}_{Q(\mathbf{X})}[(f(\mathbf{X})w(\mathbf{X}))^2] - (\mathbb{E}_{P(\mathbf{X})}[f(\mathbf{X})])^2\end{aligned}$$

# Importance Sampling (normalized)

What if you have an unnormalized distribution  $p$ ?

- Access to easy-to-sample distribution  $Q(x)$
- $Q(x) > 0$  whenever  $\tilde{P}(x) = ZP(X) > 0$
- You want to compute  $E[f(X)]$

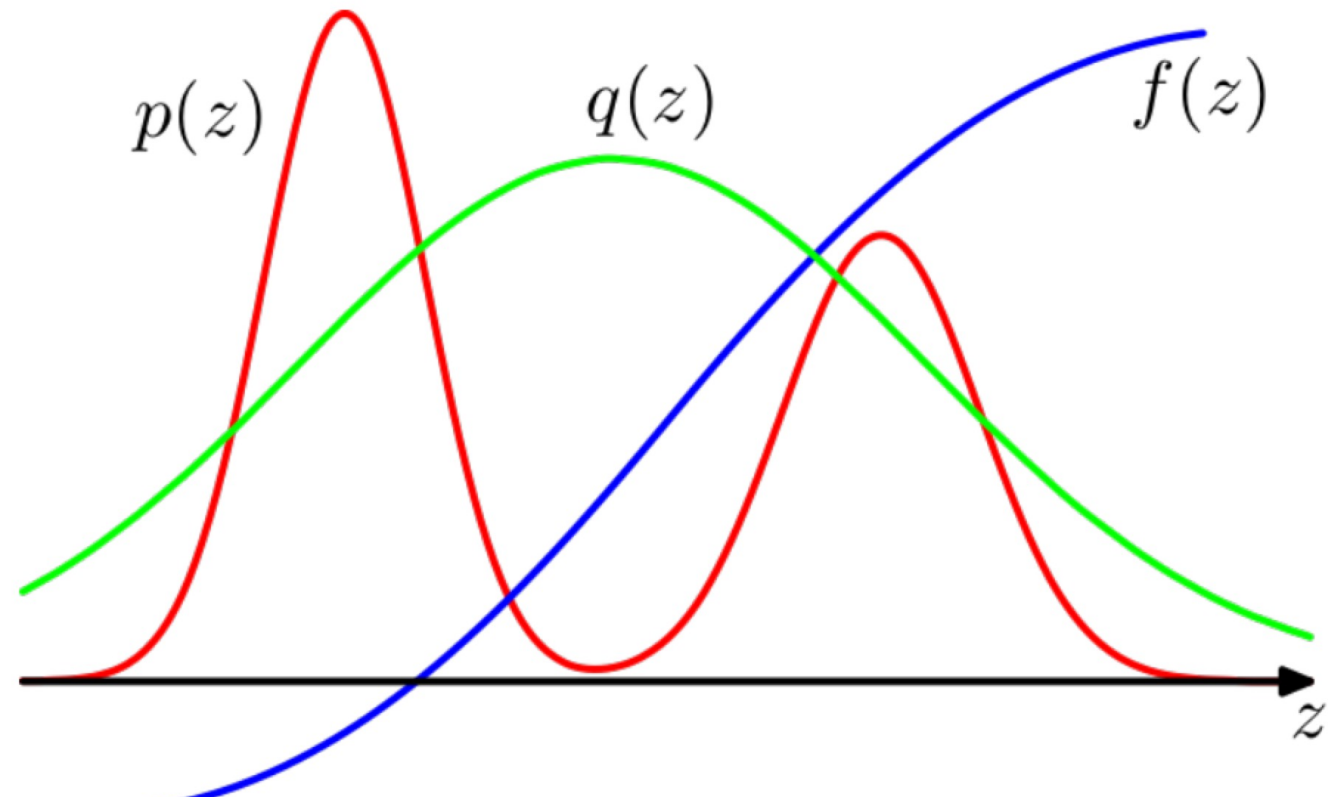


# Importance Sampling (normalized)

What if you have an unnormalized distribution  $p$ ?

- Access to easy-to-sample distribution  $Q(x)$
- $Q(x) > 0$  whenever  $\tilde{P}(x) = ZP(X) > 0$
- You want to compute  $E[f(X)]$

- $E_P[f(X)] = \frac{1}{Z} E_Q \left[ \frac{\tilde{P}(X)}{Q(X)} f(X) \right]$
- $w(x) = \frac{\tilde{P}(x)}{Q(x)}$
- $E_Q[w(X)] = Z$



# Bias/Variance of the normalized importance sampling estimator

$$\hat{\mathbb{E}}_{\mathcal{D}}(f) = \frac{\sum_{m=1}^M f(\mathbf{x}[m])w(\mathbf{x}[m])}{\sum_{m=1}^M w(\mathbf{x}[m])}$$

Not unbiased:

$$\mathbb{E}_{\mathcal{D}}[\hat{\mathbb{E}}_{\mathcal{D}}(f)] \neq \mathbb{E}_{P(\mathbf{X})}[f(\mathbf{X})]$$

Variance: 
$$\text{Var}_P[\hat{\mathbb{E}}_{\mathcal{D}}(f(\mathbf{X}))] \approx \frac{1}{M} \text{Var}_P[f(\mathbf{X})](1 + \text{Var}_Q[w(\mathbf{X})])$$

# Summary: Importance Sampling

1. Simulate from tractable distribution

$$\{x^{(m)}\}_{m=1}^M \sim Q(x)$$

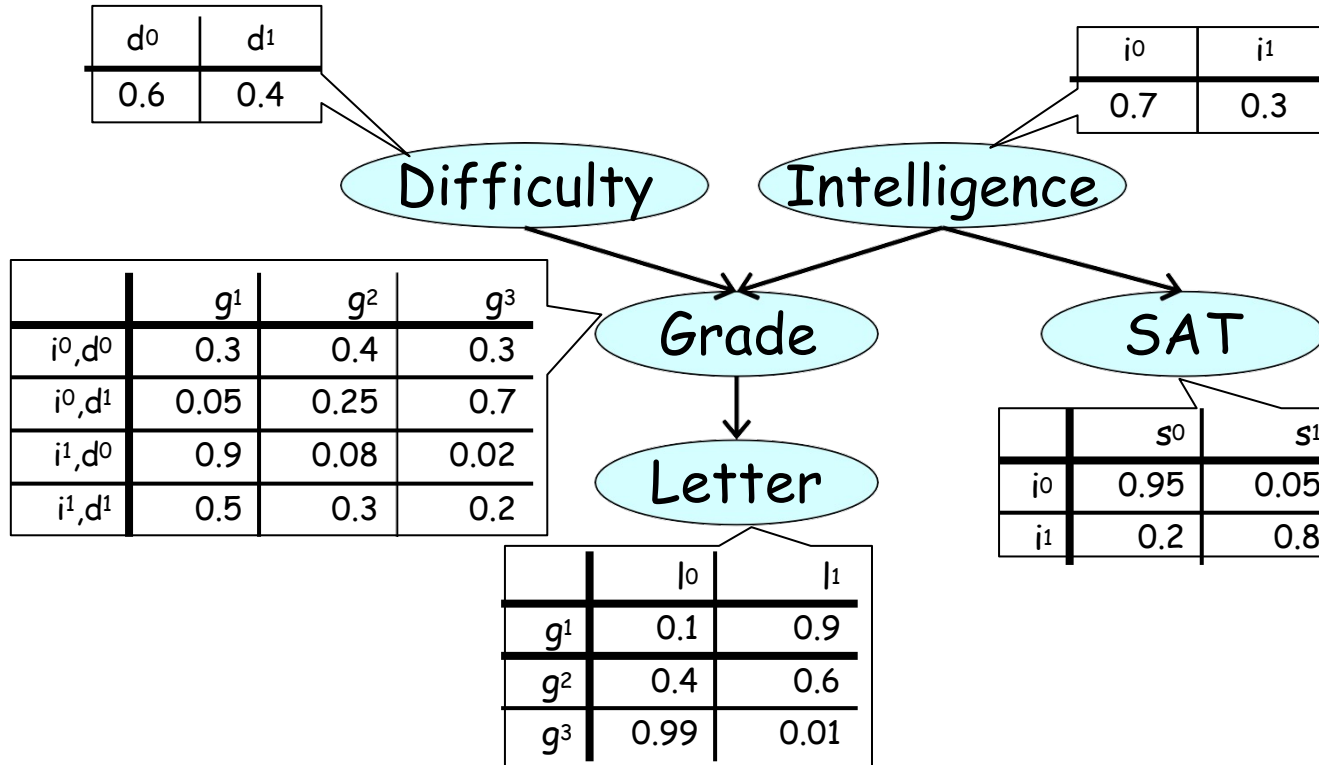
2. Compute importance weights & normalize

$$\tilde{w}^{(m)} = \frac{\tilde{P}(x^{(m)})}{Q(x^{(m)})}, w^{(m)} = \frac{\tilde{w}^{(m)}}{\sum_{i=1}^M \tilde{r}^{(i)}}$$

3. Compute importance-weighted expectation

$$\mathbf{E}_P[f(x)] \approx \sum_{m=1}^M w^{(m)} f(x^{(m)}) \equiv \hat{f}$$

# Importance Sampling in BNs



We want to sample from  $P(D, I, S, L | G = g^2)$

We want to bias our sampling towards parts of the space where this event holds

For L, this is easy; Sample from  $P(L | G = g^2)$

For the rest?

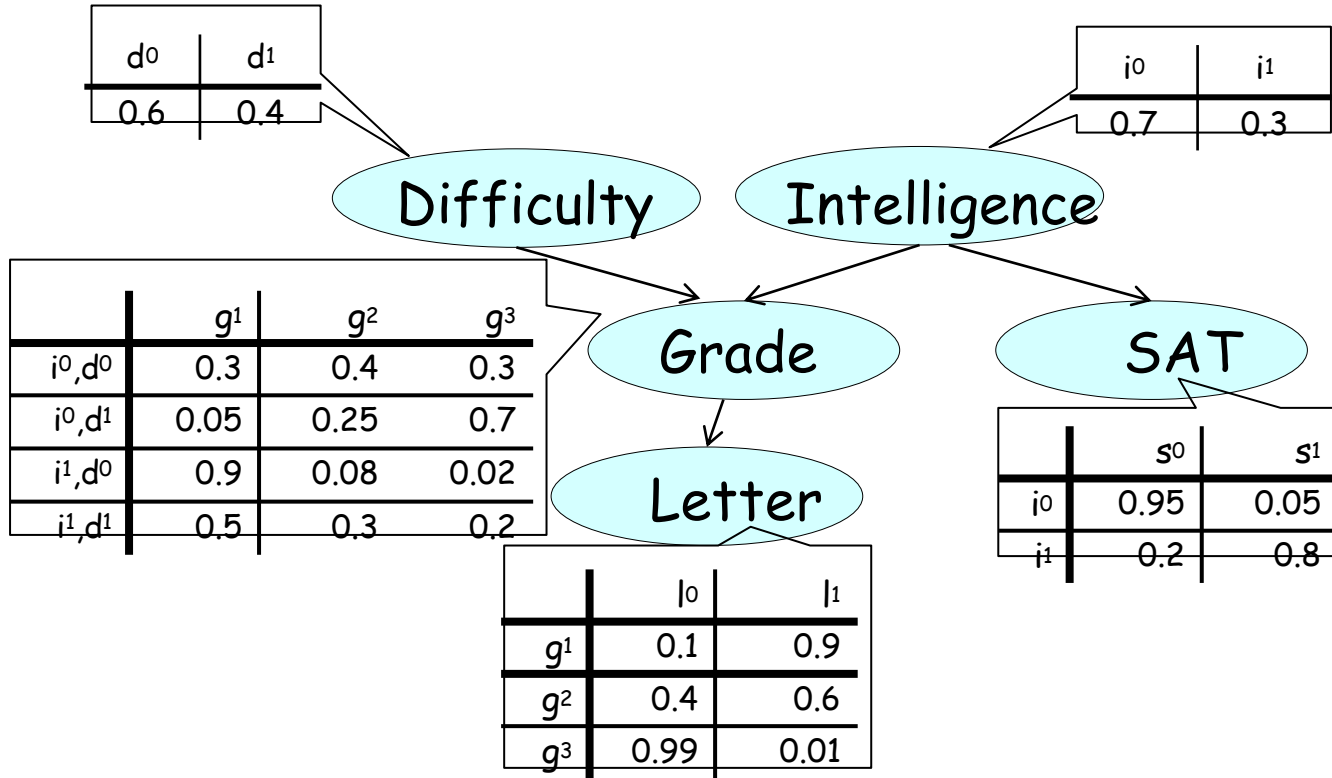
# Mutilated Network Proposal Distribution

Let  $\mathcal{B}$  be a network, and  $Z_1 = z_1, \dots, Z_k = z_k$ , abbreviated  $\mathbf{Z} = \mathbf{z}$ , an instantiation of variables.

We define the mutilated network  $\mathcal{B}_{\mathbf{Z}=\mathbf{z}}$  as follows:

- Each node  $Z_i \in \mathbf{Z}$  has no parents in  $\mathcal{B}_{\mathbf{Z}=\mathbf{z}}$ ; the CPD of  $Z_i$  in  $\mathcal{B}_{\mathbf{Z}=\mathbf{z}}$  gives probability 1 to  $Z_i = z_i$  and probability 0 to all other values  $z'_i \in \text{Val}(Z_i)$ .
- The parents and CPDs of all other nodes  $X \notin \mathbf{Z}$  are unchanged.

# Mutilated Network Proposal Distribution

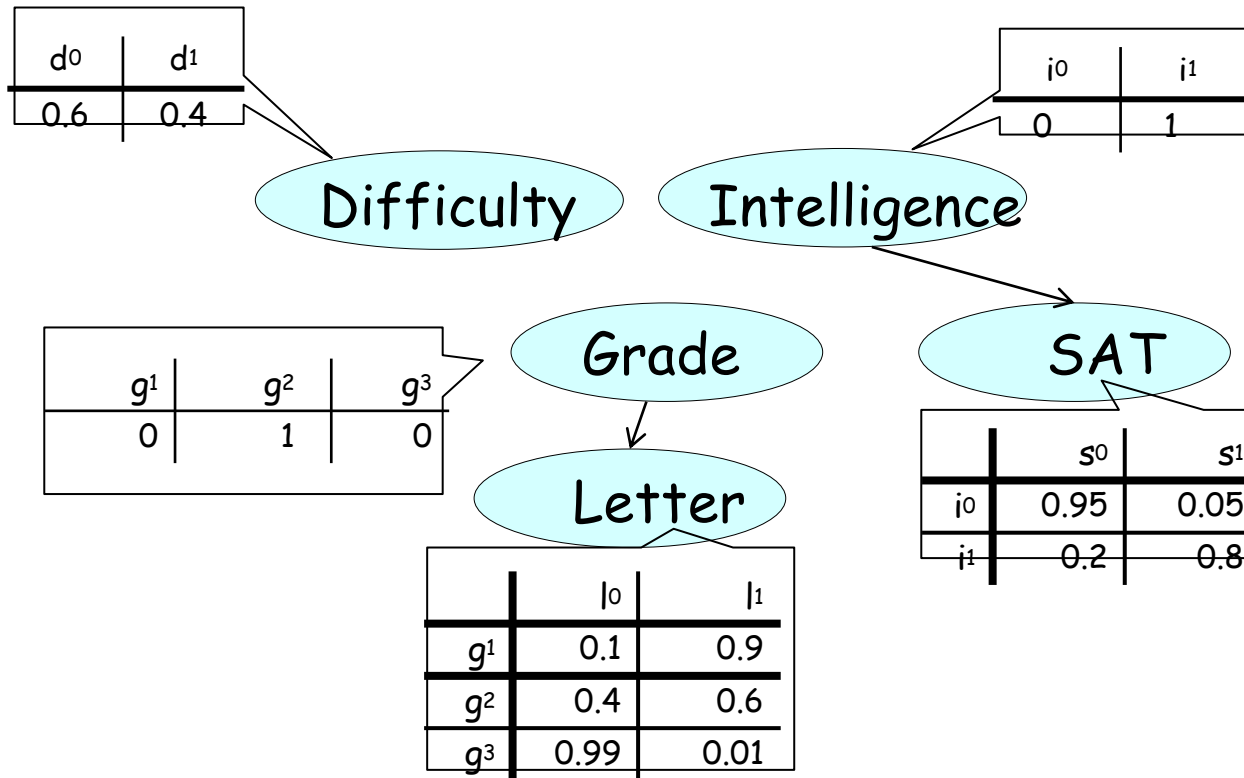


Assume you want to sample from

$$I = i^1, G = g^2$$



# Mutilated Network Proposal Distribution



Assume you want to sample from

$$I = i^1, G = g^2$$

Mutilated Graph  $\mathcal{B}_{Z=Z}$ , proposal distribution  $P_{\mathcal{B}_{Z=Z}}$

Weight of a sample  $w(\xi) = \frac{P_{\mathcal{B}}(\xi)}{P_{\mathcal{B}_{Z=Z}}(\xi)}$ .

# Summary

- Generating samples from a BN is easy
- $(\epsilon, \delta)$ -bounds exist, but usefulness is limited:
  - Additive bounds: useless for low probability events
  - Multiplicative bounds: # samples grows as  $1/P(y)$
- Forward sampling generally infeasible for MNs

Two ideas for sampling with evidence:

Throw away samples (rejection sampling)

Weigh samples (Importance sampling)

Daphne Koller

# Summary

Monte Carlo Estimators can compute unbiased expectations with good properties

Rejection sampling is straightforward but can be very wasteful

Importance sampling can reduce the variance of the estimators if you find a good proposal