

# CPSC 535

## Importance Sampling Methods

AD

8th February 2007

- Importance Sampling

- Importance Sampling
- Normalized Importance Sampling.

- Importance Sampling
- Normalized Importance Sampling.
- Importance Sampling versus Rejection Sampling.

- Let  $\pi(x)$  be a probability density on  $\mathcal{X}$ .

- Let  $\pi(x)$  be a probability density on  $\mathcal{X}$ .
- Monte Carlo approximation is given by

$$\hat{\pi}_N(x) = \frac{1}{N} \sum_{i=1}^N \delta_{X^{(i)}}(x) \text{ where } X^{(i)} \stackrel{\text{i.i.d.}}{\sim} \pi.$$

- Let  $\pi(x)$  be a probability density on  $\mathcal{X}$ .
- Monte Carlo approximation is given by

$$\hat{\pi}_N(x) = \frac{1}{N} \sum_{i=1}^N \delta_{X^{(i)}}(x) \text{ where } X^{(i)} \stackrel{\text{i.i.d.}}{\sim} \pi.$$

- For any  $\varphi : \mathcal{X} \rightarrow \mathbb{R}$

$$\mathbb{E}_{\hat{\pi}_N}(\varphi(X)) = \frac{1}{N} \sum_{i=1}^N \varphi(X^{(i)}) \approx \mathbb{E}_{\pi}(\varphi(X))$$

- Let  $\pi(x)$  be a probability density on  $\mathcal{X}$ .
- Monte Carlo approximation is given by

$$\hat{\pi}_N(x) = \frac{1}{N} \sum_{i=1}^N \delta_{X^{(i)}}(x) \text{ where } X^{(i)} \stackrel{\text{i.i.d.}}{\sim} \pi.$$

- For any  $\varphi : \mathcal{X} \rightarrow \mathbb{R}$

$$\mathbb{E}_{\hat{\pi}_N}(\varphi(X)) = \frac{1}{N} \sum_{i=1}^N \varphi(X^{(i)}) \approx \mathbb{E}_{\pi}(\varphi(X))$$

- More precisely, we have

$$\begin{aligned} \mathbb{E}_{\{X^{(i)}\}} [\mathbb{E}_{\hat{\pi}_N}(\varphi(X))] &= \mathbb{E}_{\pi}(\varphi(X)), \\ \text{var}_{\{X^{(i)}\}} (\mathbb{E}_{\hat{\pi}_N}(\varphi(X))) &= \frac{\text{var}_{\pi}(\varphi(X))}{N}. \end{aligned}$$



# Accept-Reject Procedure

- Direct methods feasible for standard distributions: inverse method, composition, etc.

# Accept-Reject Procedure

- Direct methods feasible for standard distributions: inverse method, composition, etc.
- In case where  $\pi \propto \pi^*$  does not admit any standard form, we can use a *proposal* distribution  $q$  on  $\mathcal{X}$  where  $q \propto q^*$ .

# Accept-Reject Procedure

- Direct methods feasible for standard distributions: inverse method, composition, etc.
- In case where  $\pi \propto \pi^*$  does not admit any standard form, we can use a *proposal* distribution  $q$  on  $\mathcal{X}$  where  $q \propto q^*$ .
- We need  $q$  to 'dominate'  $\pi$ ; i.e.

$$C = \sup_{x \in \mathcal{X}} \frac{\pi^*(x)}{q^*(x)} < +\infty.$$

Consider  $C' \geq C$ . Then the accept/reject procedure proceeds as follows:

### Accept/Reject procedure

- 1 Sample  $Y \sim q$  and  $U \sim \mathcal{U}(0, 1)$ .

Consider  $C' \geq C$ . Then the accept/reject procedure proceeds as follows:

### Accept/Reject procedure

- 1 Sample  $Y \sim q$  and  $U \sim \mathcal{U}(0, 1)$ .
- 2 If  $U < \frac{\pi^*(Y)}{C'q^*(Y)}$  then return  $Y$ ; otherwise return to step 1.

- This is a simple generic algorithm but it requires coming up with a bound  $C$ .

- This is a simple generic algorithm but it requires coming up with a bound  $C$ .
- Its performance typically degrade exponentially fast with the dimension of  $\mathcal{X}$ .

- This is a simple generic algorithm but it requires coming up with a bound  $C$ .
- Its performance typically degrades exponentially fast with the dimension of  $\mathcal{X}$ .
- It seems you are wasting some information by rejecting samples.



- This is a simple generic algorithm but it requires coming up with a bound  $C$ .
- Its performance typically degrades exponentially fast with the dimension of  $\mathcal{X}$ .
- It seems you are wasting some information by rejecting samples.
- You need to wait a random time to obtain some samples from  $\pi$ .

- This is a simple generic algorithm but it requires coming up with a bound  $C$ .
- Its performance typically degrades exponentially fast with the dimension of  $\mathcal{X}$ .
- It seems you are wasting some information by rejecting samples.
- You need to wait a random time to obtain some samples from  $\pi$ .
- Is it possible to “recycle” these samples?

# Importance Sampling

- Consider again the target distribution  $\pi$  and the proposal distribution  $q$ . We only require

$$\pi(x) > 0 \Rightarrow q(x) > 0.$$

# Importance Sampling

- Consider again the target distribution  $\pi$  and the proposal distribution  $q$ . We only require

$$\pi(x) > 0 \Rightarrow q(x) > 0.$$

- In this case, the Importance Sampling (IS) identity is

$$\pi(x) = w(x) q(x)$$

where the so-called Importance Weight is given by

$$w(x) = \frac{\pi(x)}{q(x)}$$

# Importance Sampling

- Consider again the target distribution  $\pi$  and the proposal distribution  $q$ . We only require

$$\pi(x) > 0 \Rightarrow q(x) > 0.$$

- In this case, the Importance Sampling (IS) identity is

$$\pi(x) = w(x) q(x)$$

where the so-called Importance Weight is given by

$$w(x) = \frac{\pi(x)}{q(x)}$$

- It follows that

$$\begin{aligned} \mathbb{E}_{\pi}(\varphi(X)) &= \int_{\mathcal{X}} \varphi(x) \pi(x) dx = \int_{\mathcal{X}} \varphi(x) \frac{\pi(x)}{q(x)} q(x) dx \\ &= \mathbb{E}_q(w(X) \varphi(X)) \end{aligned}$$

- Monte Carlo approximation of  $q$  is

$$\hat{q}_N(x) = \frac{1}{N} \sum_{i=1}^N \delta_{X^{(i)}}(x) \text{ where } X^{(i)} \stackrel{\text{i.i.d.}}{\sim} q.$$

- Monte Carlo approximation of  $q$  is

$$\hat{q}_N(x) = \frac{1}{N} \sum_{i=1}^N \delta_{X^{(i)}}(x) \text{ where } X^{(i)} \stackrel{\text{i.i.d.}}{\sim} q.$$

- It corresponds to the following approximation

$$\hat{\pi}_N(x) = \frac{1}{N} \sum_{i=1}^N w(X^{(i)}) \delta_{X^{(i)}}(x).$$

- Monte Carlo approximation of  $q$  is

$$\hat{q}_N(x) = \frac{1}{N} \sum_{i=1}^N \delta_{X^{(i)}}(x) \text{ where } X^{(i)} \stackrel{\text{i.i.d.}}{\sim} q.$$

- It corresponds to the following approximation

$$\hat{\pi}_N(x) = \frac{1}{N} \sum_{i=1}^N w(X^{(i)}) \delta_{X^{(i)}}(x).$$

- It follows that an estimate of  $\mathbb{E}_\pi(\varphi(X)) = \mathbb{E}_q(w(X)\varphi(X))$  is

$$\mathbb{E}_{\hat{q}_N}(w(X)\varphi(X)) = \frac{1}{N} \sum_{i=1}^N w(X^{(i)})\varphi(X^{(i)}).$$



- We have

$$\mathbb{E}_{\{X^{(i)}\}} [\mathbb{E}_{\hat{q}_N} (w(X)\varphi(X))] = \mathbb{E}_{\pi} (\varphi(X))$$

and

$$\begin{aligned} \text{var}_{\{X^{(i)}\}} (\mathbb{E}_{\hat{q}_N} (\varphi(X))) &= \frac{\text{var}_q (w(X)\varphi(X))}{N} \\ &= \frac{\mathbb{E}_{\pi} (w(X)\varphi^2(X)) - \mathbb{E}_{\pi}^2 (\varphi(X))}{N} \end{aligned}$$

- We have

$$\mathbb{E}_{\{X^{(i)}\}} [\mathbb{E}_{\hat{q}_N} (w(X)\varphi(X))] = \mathbb{E}_{\pi} (\varphi(X))$$

and

$$\begin{aligned} \text{var}_{\{X^{(i)}\}} (\mathbb{E}_{\hat{q}_N} (\varphi(X))) &= \frac{\text{var}_q (w(X)\varphi(X))}{N} \\ &= \frac{\mathbb{E}_{\pi} (w(X)\varphi^2(X)) - \mathbb{E}_{\pi}^2 (\varphi(X))}{N} \end{aligned}$$

- In practice, it is recommended to ensure

$$\mathbb{E}_{\pi} (w(X)) = \int \frac{\pi^2(x)}{q(x)} dx < \infty.$$

- We have

$$\mathbb{E}_{\{X^{(i)}\}} [\mathbb{E}_{\hat{q}_N} (w(X)\varphi(X))] = \mathbb{E}_{\pi} (\varphi(X))$$

and

$$\begin{aligned} \text{var}_{\{X^{(i)}\}} (\mathbb{E}_{\hat{q}_N} (\varphi(X))) &= \frac{\text{var}_q (w(X)\varphi(X))}{N} \\ &= \frac{\mathbb{E}_{\pi} (w(X)\varphi^2(X)) - \mathbb{E}_{\pi}^2 (\varphi(X))}{N} \end{aligned}$$

- In practice, it is recommended to ensure

$$\mathbb{E}_{\pi} (w(X)) = \int \frac{\pi^2(x)}{q(x)} dx < \infty.$$

- Even if it is not necessary, it is actually even better to ensure that

$$\sup_{x \in \mathcal{X}} w(x) < \infty.$$

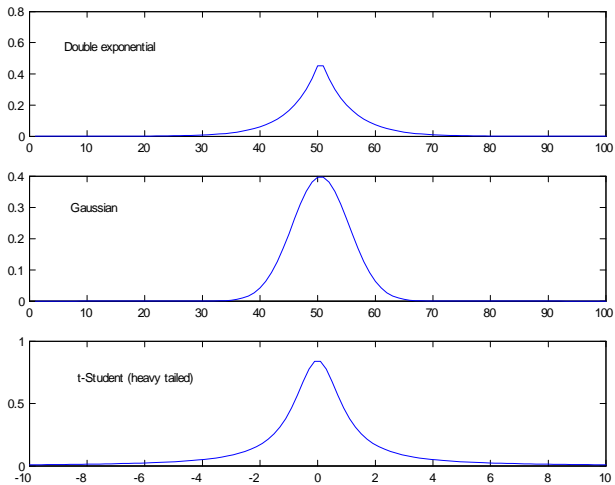


Figure: Target double exponential distributions and two IS distributions

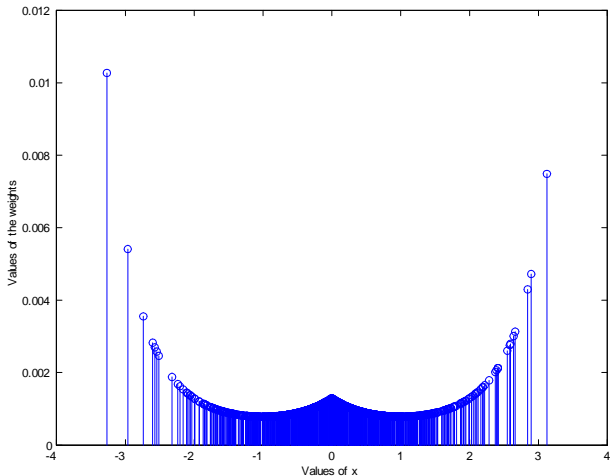


Figure: IS approximation obtained using a Gaussian IS distribution

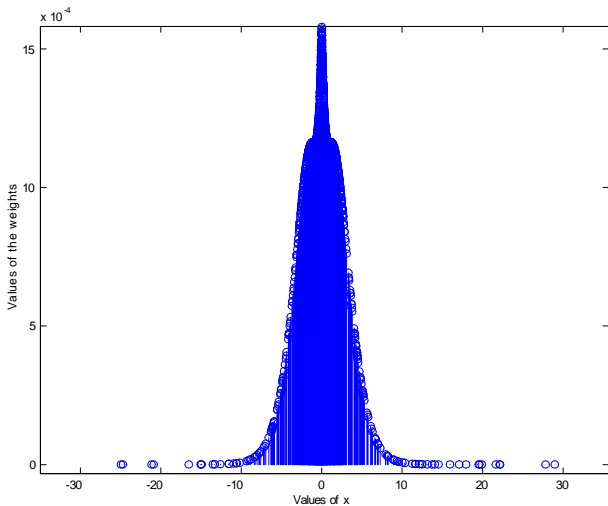


Figure: IS approximation obtained using a Student-t IS distribution

- We try to compute

$$\int \sqrt{\frac{x}{1-x}} \pi(x) dx$$

where

$$\pi(x) = \frac{\Gamma((\nu+1)/2)}{\sqrt{\nu\pi}\Gamma(\nu/2)} \left(1 + \frac{x}{\nu}\right)^{-(\nu+1)/2}$$

is a t-student distribution with  $\nu > 1$  (you can sample from it by composition  $\mathcal{N}(0, 1) / \mathcal{G}a(\nu/2, \nu/2)$ ) using Monte Carlo.

- We try to compute

$$\int \sqrt{\frac{x}{1-x}} \pi(x) dx$$

where

$$\pi(x) = \frac{\Gamma((\nu+1)/2)}{\sqrt{\nu\pi}\Gamma(\nu/2)} \left(1 + \frac{x}{\nu}\right)^{-(\nu+1)/2}$$

is a t-student distribution with  $\nu > 1$  (you can sample from it by composition  $\mathcal{N}(0, 1) / \mathcal{G}a(\nu/2, \nu/2)$ ) using Monte Carlo.

- We use  $q_1(x) = \pi(x)$ ,  $q_2(x) = \frac{\Gamma(1)}{\sqrt{\nu\pi}\Gamma(1/2)} \left(1 + \frac{x}{\nu}\right)^{-1}$  (Cauchy distribution) and  $q_3(x) = \mathcal{N}(x; 0, \frac{\nu}{\nu-2})$  (variance chosen to match the variance of  $\pi(x)$ )



- We try to compute

$$\int \sqrt{\frac{x}{1-x}} \pi(x) dx$$

where

$$\pi(x) = \frac{\Gamma((\nu+1)/2)}{\sqrt{\nu\pi}\Gamma(\nu/2)} \left(1 + \frac{x}{\nu}\right)^{-(\nu+1)/2}$$

is a t-student distribution with  $\nu > 1$  (you can sample from it by composition  $\mathcal{N}(0, 1) / \mathcal{G}a(\nu/2, \nu/2)$ ) using Monte Carlo.

- We use  $q_1(x) = \pi(x)$ ,  $q_2(x) = \frac{\Gamma(1)}{\sqrt{\nu\pi}\Gamma(1/2)} \left(1 + \frac{x}{\nu\sigma}\right)^{-1}$  (Cauchy distribution) and  $q_3(x) = \mathcal{N}(x; 0, \frac{\nu}{\nu-2})$  (variance chosen to match the variance of  $\pi(x)$ )
- It is easy to see that

$$\frac{\pi(x)}{q_2(x)} < \infty \text{ and } \int \frac{\pi(x)^2}{q_3(x)} dx = \infty, \quad \frac{\pi(x)}{q_3(x)} \text{ is unbounded}$$

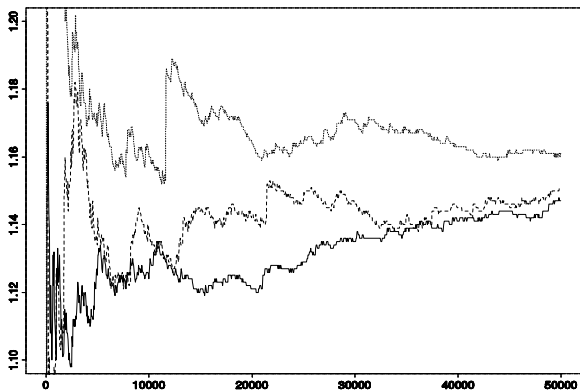


Figure: Performance for  $\nu = 12$  with  $q_1$  (solid line),  $q_2$  (dashes) and  $q_3$  (light dots). Final values 1.14, 1.14 and 1.16 vs true value 1.13

- We now try to compute

$$\int_{2.1}^{\infty} x^5 \pi(x) dx$$

- We now try to compute

$$\int_{2.1}^{\infty} x^5 \pi(x) dx$$

- We try to use the same importance distribution but also use the fact that using a change of variables  $u = 1/x$ , we have

$$\begin{aligned} \int_{2.1}^{\infty} x^5 \pi(x) dx &= \int_0^{1/2.1} u^{-7} \pi(1/u) du \\ &= \frac{1}{2.1} \int_0^{1/2.1} 2.1 u^{-7} \pi(1/u) du \end{aligned}$$

which is the expectation of  $2.1 u^{-7} \pi(1/u)$  with respect to  $\mathcal{U}[0, 1/2.1]$ .

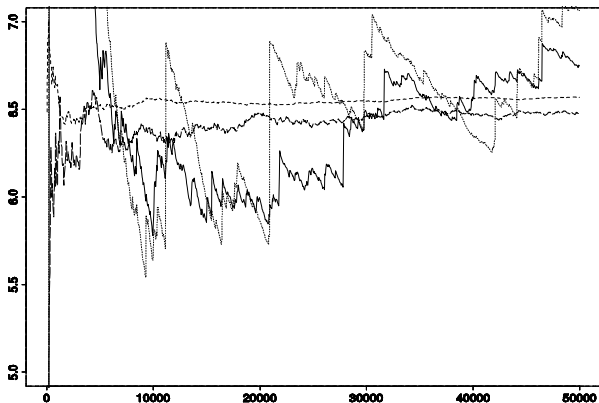


Figure: Performance for  $\nu = 12$  with  $q_1$  (solid),  $q_2$  (short dashes),  $q_3$  (dots), uniform (long dashes). Final values 6.75, 6.48, 7.06 and 6.48 vs true value 6.54

# Optimal Importance function

- For a given test function, one can minimize the IS variance using

$$q^{\text{opt}}(x) = \frac{|\varphi(x)| \pi(x)}{\int_{\mathcal{X}} |\varphi(x)| \pi(x) dx}$$

# Optimal Importance function

- For a given test function, one can minimize the IS variance using

$$q^{\text{opt}}(x) = \frac{|\varphi(x)| \pi(x)}{\int_{\mathcal{X}} |\varphi(x)| \pi(x) dx}$$

- *Proof.*

$$\text{var}_q(w(X)\varphi(X)) = \int q(x) \frac{\pi^2(x)}{q^2(x)} \varphi^2(x) dx - \left( \int \pi(x) \varphi(x) dx \right)^2$$

and

$$\begin{aligned} \int q(x) \frac{\pi^2(x)}{q^2(x)} \varphi^2(x) dx &\geq \left( \int q(x) \frac{\pi(x) |\varphi(x)|}{q(x)} dx \right)^2 \\ &= \left( \int \pi(x) |\varphi(x)| dx \right)^2. \end{aligned}$$

This lower bound is attained for  $q^{\text{opt}}(x)$ .

# Normalized Importance Sampling

- In most if not all applications we are interested in, standard IS cannot be used as the importance weights  $w(x) = \pi(x) / q(x)$  cannot be evaluated in closed-form. In practice, we typically only know  $\pi(x) \propto \pi^*(x)$  and  $q(x) \propto q^*(x)$ .



# Normalized Importance Sampling

- In most if not all applications we are interested in, standard IS cannot be used as the importance weights  $w(x) = \pi(x) / q(x)$  cannot be evaluated in closed-form. In practice, we typically only know  $\pi(x) \propto \pi^*(x)$  and  $q(x) \propto q^*(x)$ .
- **Normalized IS** identity is based on

$$\begin{aligned}\pi(x) &= \frac{\pi^*(x)}{\int \pi^*(x) dx} = \frac{w^*(x) q^*(x)}{\int w^*(x) q^*(x) dx} \\ &= \frac{w^*(x) q(x)}{\int w^*(x) q(x) dx} = \frac{w(x) q(x)}{\int w(x) q(x) dx}\end{aligned}$$

where

$$w^*(x) = \frac{\pi^*(x)}{q^*(x)}.$$

- For any test function  $\varphi$ , we can also write

$$\mathbb{E}_{\pi}(\varphi(X)) = \frac{\mathbb{E}_q(w^*(X)\varphi(X))}{\mathbb{E}_q(w^*(X))} = \frac{\mathbb{E}_q(w(X)\varphi(X))}{\mathbb{E}_q(w(X))}.$$

- For any test function  $\varphi$ , we can also write

$$\mathbb{E}_{\pi}(\varphi(X)) = \frac{\mathbb{E}_q(w^*(X)\varphi(X))}{\mathbb{E}_q(w^*(X))} = \frac{\mathbb{E}_q(w(X)\varphi(X))}{\mathbb{E}_q(w(X))}.$$

- Given a Monte Carlo approximation of  $q$ ;  $\hat{q}_N(x) = \frac{1}{N} \sum_{i=1}^N \delta_{X^{(i)}}(x)$  then

$$\hat{\pi}_N(x) = \sum_{i=1}^N W^{(i)} \delta_{X^{(i)}}(x) \text{ where } W^{(i)} = \frac{w^*(X^{(i)})}{\sum_{j=1}^N w^*(X^{(j)})},$$

$$\mathbb{E}_{\hat{\pi}_N}(\varphi(X)) = \sum_{i=1}^N W^{(i)} \varphi(X^{(i)}).$$

- For any test function  $\varphi$ , we can also write

$$\mathbb{E}_{\pi}(\varphi(X)) = \frac{\mathbb{E}_q(w^*(X)\varphi(X))}{\mathbb{E}_q(w^*(X))} = \frac{\mathbb{E}_q(w(X)\varphi(X))}{\mathbb{E}_q(w(X))}.$$

- Given a Monte Carlo approximation of  $q$ ;  $\hat{q}_N(x) = \frac{1}{N} \sum_{i=1}^N \delta_{X^{(i)}}(x)$  then

$$\hat{\pi}_N(x) = \sum_{i=1}^N W^{(i)} \delta_{X^{(i)}}(x) \text{ where } W^{(i)} = \frac{w^*(X^{(i)})}{\sum_{j=1}^N w^*(X^{(j)})},$$

$$\mathbb{E}_{\hat{\pi}_N}(\varphi(X)) = \sum_{i=1}^N W^{(i)} \varphi(X^{(i)}).$$

- The estimates are a ratio of estimates.

- Contrary to standard IS, this estimate is biased but by the LLN it is asymptotically consistent.

- Contrary to standard IS, this estimate is biased but by the LLN it is asymptotically consistent.
- Derivation of the asymptotic bias and variance based on the delta method.

# Asymptotic Bias and Variance

- Assume you have  $Z = g(A, B)$  with  $\mathbb{E}(A) = \mu_A$  and  $\mathbb{E}(B) = \mu_B$  then a two-dimensional Taylor expansion gives around  $\mu = (\mu_A, \mu_B)$

$$Z \approx g(\mu) + (A - \mu_A) \frac{\partial g}{\partial a}(\mu) + (B - \mu_B) \frac{\partial g}{\partial b}(\mu).$$

It follows that

$$\mathbb{E}(Z) \approx g(\mu),$$

$$\text{Var}(Z) \approx \sigma_A^2 \frac{\partial g^2}{\partial a}(\mu) + \sigma_B^2 \frac{\partial g^2}{\partial b}(\mu) + 2 \frac{\partial g}{\partial a}(\mu) \frac{\partial g}{\partial b}(\mu) \sigma_{A,B}.$$

# Asymptotic Bias and Variance

- Assume you have  $Z = g(A, B)$  with  $\mathbb{E}(A) = \mu_A$  and  $\mathbb{E}(B) = \mu_B$  then a two-dimensional Taylor expansion gives around  $\mu = (\mu_A, \mu_B)$

$$Z \approx g(\mu) + (A - \mu_A) \frac{\partial g}{\partial a}(\mu) + (B - \mu_B) \frac{\partial g}{\partial b}(\mu).$$

It follows that

$$\mathbb{E}(Z) \approx g(\mu),$$

$$\text{Var}(Z) \approx \sigma_A^2 \frac{\partial g^2}{\partial a}(\mu) + \sigma_B^2 \frac{\partial g^2}{\partial b}(\mu) + 2 \frac{\partial g}{\partial a}(\mu) \frac{\partial g}{\partial b}(\mu) \sigma_{A,B}.$$

- In our case

$$Z = \mathbb{E}_{\hat{\pi}_N}(\varphi(X)) = \frac{\mathbb{E}_{\hat{q}_N}(w^*(X) \varphi(X))}{\mathbb{E}_{\hat{q}_N}(w^*(X))} = \frac{A}{B}$$



- We have

$$\frac{\partial g}{\partial a}(\mu) \frac{\partial g}{\partial b}(\mu) = -\frac{\mu_A}{\mu_B^3}, \quad \frac{\partial g^2}{\partial a}(\mu) = \frac{1}{\mu_B^2}, \quad \frac{\partial g^2}{\partial b}(\mu) = \frac{\mu_A^2}{\mu_B^4},$$

$$\begin{aligned} \mu_A &= \mathbb{E}_q(w^*(X) \varphi(X)), \quad \mu_B = \mathbb{E}_q(w^*(X)), \\ \sigma_A^2 &= \frac{\text{var}_q(w^*(X) \varphi(X))}{N}, \quad \sigma_B^2 = \frac{\text{var}_q(w^*(X))}{N} \end{aligned}$$

$$\sigma_{A,B} = \frac{\mathbb{E}_q(w^*(X)^2 \varphi(X)) - \mu_A \cdot \mu_B}{N}.$$

- It follows that

$$\begin{aligned} & \text{var} (\mathbb{E}_{\hat{\pi}_N} (\varphi (X))) \\ & \approx \sigma_A^2 \frac{\partial g^2}{\partial a} (\mu) + \sigma_B^2 \frac{\partial g^2}{\partial b} (\mu) + 2 \frac{\partial g}{\partial a} (\mu) \frac{\partial g}{\partial b} (\mu) \sigma_{A,B} \\ & = \frac{\sigma_A^2}{\mu_B^2} + \frac{\sigma_B^2 \mu_A^2}{\mu_B^4} - 2 \frac{\mu_A \sigma_{A,B}}{\mu_B^3} \end{aligned}$$

- It follows that

$$\begin{aligned}
 & \text{var} (\mathbb{E}_{\hat{\pi}_N} (\varphi (X))) \\
 & \approx \sigma_A^2 \frac{\partial g^2}{\partial a} (\mu) + \sigma_B^2 \frac{\partial g^2}{\partial b} (\mu) + 2 \frac{\partial g}{\partial a} (\mu) \frac{\partial g}{\partial b} (\mu) \sigma_{A,B} \\
 & = \frac{\sigma_A^2}{\mu_B^2} + \frac{\sigma_B^2 \mu_A^2}{\mu_B^4} - 2 \frac{\mu_A \sigma_{A,B}}{\mu_B^3}
 \end{aligned}$$

- Asymptotically, we have a central limit theorem

$$\sqrt{N} (\mathbb{E}_{\hat{\pi}_N} (\varphi (X)) - \mathbb{E}_{\pi} (\varphi (X))) \Rightarrow \mathcal{N} (0, \sigma_{IS}^2 (\varphi))$$

where

$$\sigma_{IS}^2 (\varphi) = \int \frac{\pi^2 (x)}{q (x)} (\varphi (x) - \mathbb{E}_{\pi} (\varphi))^2 dx$$

- In practice, it is now necessary but highly recommended to select the proposal  $q$  such that

$$\sup_{x \in \mathcal{X}} w(x) < \infty \text{ or equivalently } \sup_{x \in \mathcal{X}} w^*(x) < \infty.$$

- In practice, it is now necessary but highly recommended to select the proposal  $q$  such that

$$\sup_{x \in \mathcal{X}} w(x) < \infty \text{ or equivalently } \sup_{x \in \mathcal{X}} w^*(x) < \infty.$$

- There is some empirical evidence that Normalized IS performs better than standard IS in numerous cases.

- Using a second order Taylor expansion

$$\begin{aligned} Z &\approx g(\mu) + (A - \mu_A) \frac{\partial g}{\partial a}(\mu) + (B - \mu_B) \frac{\partial g}{\partial b}(\mu) \\ &+ \frac{1}{2} (A - \mu_A)^2 \frac{\partial^2 g}{\partial a^2}(\mu) + \frac{1}{2} (B - \mu_B)^2 \frac{\partial^2 g}{\partial b^2}(\mu) \\ &+ (A - \mu_A)(B - \mu_B) \frac{\partial^2 g}{\partial a \partial b}(\mu) \end{aligned}$$

gives

$$\begin{aligned} \mathbb{E}(\mathbb{E}_{\hat{\pi}_N}(\varphi(X))) &\approx g(\mu) + \frac{1}{2} \sigma_A^2 \frac{\partial^2 g}{\partial a^2}(\mu) + \frac{1}{2} \sigma_B^2 \frac{\partial^2 g}{\partial b^2}(\mu) \\ &+ \sigma_{A,B} \frac{\partial^2 g}{\partial a \partial b}(\mu). \end{aligned}$$

- Using a second order Taylor expansion

$$\begin{aligned} Z &\approx g(\mu) + (A - \mu_A) \frac{\partial g}{\partial a}(\mu) + (B - \mu_B) \frac{\partial g}{\partial b}(\mu) \\ &+ \frac{1}{2} (A - \mu_A)^2 \frac{\partial^2 g}{\partial a^2}(\mu) + \frac{1}{2} (B - \mu_B)^2 \frac{\partial^2 g}{\partial b^2}(\mu) \\ &+ (A - \mu_A)(B - \mu_B) \frac{\partial^2 g}{\partial a \partial b}(\mu) \end{aligned}$$

gives

$$\begin{aligned} \mathbb{E}(\mathbb{E}_{\hat{\pi}_N}(\varphi(X))) &\approx g(\mu) + \frac{1}{2} \sigma_A^2 \frac{\partial^2 g}{\partial a^2}(\mu) + \frac{1}{2} \sigma_B^2 \frac{\partial^2 g}{\partial b^2}(\mu) \\ &+ \sigma_{A,B} \frac{\partial^2 g}{\partial a \partial b}(\mu). \end{aligned}$$

- It follows that asymptotically we have

$$N(\mathbb{E}_{\hat{\pi}_N}(\varphi(X)) - \mathbb{E}_{\pi}(\varphi(X))) \rightarrow - \int \frac{\pi^2(x)}{q(x)} (\varphi(x) - \mathbb{E}_{\pi}(\varphi)) dx.$$

- Using a second order Taylor expansion

$$\begin{aligned} Z &\approx g(\mu) + (A - \mu_A) \frac{\partial g}{\partial a}(\mu) + (B - \mu_B) \frac{\partial g}{\partial b}(\mu) \\ &+ \frac{1}{2} (A - \mu_A)^2 \frac{\partial^2 g}{\partial a^2}(\mu) + \frac{1}{2} (B - \mu_B)^2 \frac{\partial^2 g}{\partial b^2}(\mu) \\ &+ (A - \mu_A)(B - \mu_B) \frac{\partial^2 g}{\partial a \partial b}(\mu) \end{aligned}$$

gives

$$\begin{aligned} \mathbb{E}(\mathbb{E}_{\hat{\pi}_N}(\varphi(X))) &\approx g(\mu) + \frac{1}{2} \sigma_A^2 \frac{\partial^2 g}{\partial a^2}(\mu) + \frac{1}{2} \sigma_B^2 \frac{\partial^2 g}{\partial b^2}(\mu) \\ &+ \sigma_{A,B} \frac{\partial^2 g}{\partial a \partial b}(\mu). \end{aligned}$$

- It follows that asymptotically we have

$$N(\mathbb{E}_{\hat{\pi}_N}(\varphi(X)) - \mathbb{E}_{\pi}(\varphi(X))) \rightarrow - \int \frac{\pi^2(x)}{q(x)} (\varphi(x) - \mathbb{E}_{\pi}(\varphi)) dx.$$

- We have  $Bias^2$  of order  $1/N^2$  and Variance of order  $1/N$ .



- The asymptotic variance (and also the asymptotic bias) can be consistently estimated from the data using

$$\frac{\widehat{\sigma_{IS}^2}(\varphi)}{N} = \frac{\widehat{\sigma}_A^2}{\widehat{\mu}_B^2} + \frac{\widehat{\sigma}_B^2 \widehat{\mu}_A^2}{\widehat{\mu}_B^4} - 2 \frac{\widehat{\mu}_A \widehat{\sigma}_{A,B}}{\widehat{\mu}_B^3}.$$

- The asymptotic variance (and also the asymptotic bias) can be consistently estimated from the data using

$$\frac{\widehat{\sigma_{IS}^2}(\varphi)}{N} = \frac{\widehat{\sigma}_A^2}{\widehat{\mu}_B^2} + \frac{\widehat{\sigma}_B^2 \widehat{\mu}_A^2}{\widehat{\mu}_B^4} - 2 \frac{\widehat{\mu}_A \widehat{\sigma}_{A,B}}{\widehat{\mu}_B^3}.$$

- You can also compute the variance of the variance estimate; see Geweke (1989).

# Application to Bayesian Statistics

- Consider a Bayesian model: prior  $\pi(\theta)$  and likelihood  $f(x|\theta)$ .

# Application to Bayesian Statistics

- Consider a Bayesian model: prior  $\pi(\theta)$  and likelihood  $f(x|\theta)$ .
- The posterior distribution is given by

$$\pi(\theta|x) = \frac{\pi(\theta)f(x|\theta)}{\int_{\Theta} \pi(\theta)f(x|\theta)d\theta} \propto \pi^*(\theta|x)$$

where  $\pi^*(\theta|x) = \pi(\theta)f(x|\theta)$ .

# Application to Bayesian Statistics

- Consider a Bayesian model: prior  $\pi(\theta)$  and likelihood  $f(x|\theta)$ .
- The posterior distribution is given by

$$\pi(\theta|x) = \frac{\pi(\theta)f(x|\theta)}{\int_{\Theta} \pi(\theta)f(x|\theta)d\theta} \propto \pi^*(\theta|x)$$

where  $\pi^*(\theta|x) = \pi(\theta)f(x|\theta)$ .

- We can use the prior distribution as a candidate distribution  $q(\theta) = q^*(\theta) = \pi(\theta)$ .

# Application to Bayesian Statistics

- Consider a Bayesian model: prior  $\pi(\theta)$  and likelihood  $f(x|\theta)$ .
- The posterior distribution is given by

$$\pi(\theta|x) = \frac{\pi(\theta)f(x|\theta)}{\int_{\Theta} \pi(\theta)f(x|\theta)d\theta} \propto \pi^*(\theta|x)$$

where  $\pi^*(\theta|x) = \pi(\theta)f(x|\theta)$ .

- We can use the prior distribution as a candidate distribution  $q(\theta) = q^*(\theta) = \pi(\theta)$ .
- We also get an estimate of the marginal likelihood

$$\int_{\Theta} \pi(\theta) f(x|\theta) d\theta.$$

- *Example:* Application to Bayesian analysis of Markov chain. Consider a two state Markov chain with transition matrix  $F$

$$\begin{pmatrix} p_1 & 1 - p_1 \\ 1 - p_2 & p_2 \end{pmatrix}$$

that is  $\Pr(X_{t+1} = 1 | X_t = 1) = 1 - \Pr(X_{t+1} = 2 | X_t = 1) = p_1$  and  $\Pr(X_{t+1} = 2 | X_t = 2) = 1 - \Pr(X_{t+1} = 1 | X_t = 2) = p_2$ . Physical constraints tell us that  $p_1 + p_2 < 1$ .

- *Example:* Application to Bayesian analysis of Markov chain. Consider a two state Markov chain with transition matrix  $F$

$$\begin{pmatrix} p_1 & 1 - p_1 \\ 1 - p_2 & p_2 \end{pmatrix}$$

that is  $\Pr(X_{t+1} = 1 | X_t = 1) = 1 - \Pr(X_{t+1} = 2 | X_t = 1) = p_1$  and  $\Pr(X_{t+1} = 2 | X_t = 2) = 1 - \Pr(X_{t+1} = 1 | X_t = 2) = p_2$ . Physical constraints tell us that  $p_1 + p_2 < 1$ .

- Assume we observe  $x_1, \dots, x_m$  and the prior is

$$\pi(p_1, p_2) = 2\mathbb{I}_{p_1+p_2 \leq 1}$$

then the posterior is

$$\pi(p_1, p_2 | x_{1:m}) \propto p_1^{m_{1,1}} (1 - p_1)^{m_{1,2}} (1 - p_2)^{m_{2,1}} p_2^{m_{2,2}} \mathbb{I}_{p_1+p_2 \leq 1}$$

where

$$m_{i,j} = \sum_{t=1}^{m-1} \mathbb{I}_{x_t=i} \mathbb{I}_{x_{t+1}=j}$$



- *Example:* Application to Bayesian analysis of Markov chain. Consider a two state Markov chain with transition matrix  $F$

$$\begin{pmatrix} p_1 & 1 - p_1 \\ 1 - p_2 & p_2 \end{pmatrix}$$

that is  $\Pr(X_{t+1} = 1 | X_t = 1) = 1 - \Pr(X_{t+1} = 2 | X_t = 1) = p_1$  and  $\Pr(X_{t+1} = 2 | X_t = 2) = 1 - \Pr(X_{t+1} = 1 | X_t = 2) = p_2$ . Physical constraints tell us that  $p_1 + p_2 < 1$ .

- Assume we observe  $x_1, \dots, x_m$  and the prior is

$$\pi(p_1, p_2) = 2\mathbb{I}_{p_1+p_2 \leq 1}$$

then the posterior is

$$\pi(p_1, p_2 | x_{1:m}) \propto p_1^{m_{1,1}} (1 - p_1)^{m_{1,2}} (1 - p_2)^{m_{2,1}} p_2^{m_{2,2}} \mathbb{I}_{p_1+p_2 \leq 1}$$

where

$$m_{i,j} = \sum_{t=1}^{m-1} \mathbb{I}_{x_t=i} \mathbb{I}_{x_{t+1}=j}$$

- The posterior does not admit a standard expression and its normalizing constant is unknown. We can sample from it using rejection sampling.

- We are interested in estimating  $\mathbb{E} [\varphi_i (p_1, p_2) | x_{1:m}]$  for  
 $\varphi_1 (p_1, p_2) = p_1$ ,  $\varphi_2 (p_1, p_2) = p_2$ ,  $\varphi_3 (p_1, p_2) = p_1 / (1 - p_1)$ ,  
 $\varphi_4 (p_1, p_2) = p_2 / (1 - p_2)$  and  $\varphi_5 (p_1, p_2) = \log \frac{p_1(1-p_2)}{p_2(1-p_1)}$  using  
Importance Sampling.

- We are interested in estimating  $\mathbb{E} [\varphi_i (p_1, p_2) | x_{1:m}]$  for  $\varphi_1 (p_1, p_2) = p_1$ ,  $\varphi_2 (p_1, p_2) = p_2$ ,  $\varphi_3 (p_1, p_2) = p_1 / (1 - p_1)$ ,  $\varphi_4 (p_1, p_2) = p_2 / (1 - p_2)$  and  $\varphi_5 (p_1, p_2) = \log \frac{p_1(1-p_2)}{p_2(1-p_1)}$  using Importance Sampling.
- If there was no on  $p_1 + p_2 < 1$  and  $\pi (p_1, p_2)$  was uniform on  $[0, 1] \times [0, 1]$ , then the posterior would be

$$\begin{aligned} \pi_0 (p_1, p_2 | x_{1:m}) &= \text{Be} (p_1; m_{1,1} + 1, m_{1,2} + 1) \\ &\quad \text{Be} (p_2; m_{2,2} + 1, m_{2,1} + 1) \end{aligned}$$

but this is inefficient as for the given data  $(m_{1,1}, m_{1,2}, m_{2,2}, m_{2,1})$  we have  $\pi_0 (p_1 + p_2 < 1 | x_{1:m}) = 0.21$ .

- We are interested in estimating  $\mathbb{E} [\varphi_i (p_1, p_2) | x_{1:m}]$  for  $\varphi_1 (p_1, p_2) = p_1$ ,  $\varphi_2 (p_1, p_2) = p_2$ ,  $\varphi_3 (p_1, p_2) = p_1 / (1 - p_1)$ ,  $\varphi_4 (p_1, p_2) = p_2 / (1 - p_2)$  and  $\varphi_5 (p_1, p_2) = \log \frac{p_1(1-p_2)}{p_2(1-p_1)}$  using Importance Sampling.
- If there was no on  $p_1 + p_2 < 1$  and  $\pi (p_1, p_2)$  was uniform on  $[0, 1] \times [0, 1]$ , then the posterior would be

$$\begin{aligned} \pi_0 (p_1, p_2 | x_{1:m}) &= \mathcal{B}e (p_1; m_{1,1} + 1, m_{1,2} + 1) \\ &\quad \mathcal{B}e (p_2; m_{2,2} + 1, m_{2,1} + 1) \end{aligned}$$

but this is inefficient as for the given data  $(m_{1,1}, m_{1,2}, m_{2,2}, m_{2,1})$  we have  $\pi_0 (p_1 + p_2 < 1 | x_{1:m}) = 0.21$ .

- The form of the posterior suggests using a Dirichlet distribution with density

$$\pi_1 (p_1, p_2 | x_{1:m}) \propto p_1^{m_{1,1}} p_2^{m_{2,2}} (1 - p_1 - p_2)^{m_{1,2} + m_{2,1}}$$

but  $\pi (p_1, p_2 | x_{1:m}) / \pi_1 (p_1, p_2 | x_{1:m})$  is unbounded.

- (Geweke, 1989) proposed using the normal approximation to the binomial distribution

$$\begin{aligned} \pi_2(p_1, p_2 | x_{1:m}) &\propto \exp\left(- (m_{1,1} + m_{1,2}) (p_1 - \hat{p}_1)^2 / (2\hat{p}_1 (1 - \hat{p}_1))\right) \\ &\quad \times \exp\left(- (m_{2,1} + m_{2,2}) (p_2 - \hat{p}_2)^2 / (2\hat{p}_2 (1 - \hat{p}_2))\right) \end{aligned}$$

where  $\hat{p}_1 = m_{1,1} / (m_{1,1} + m_{1,2})$ ,  $\hat{p}_2 = m_{2,2} / (m_{2,2} + m_{2,1})$ . Then to simulate from this distribution, we simulate first  $\pi_2(p_1 | x_{1:m})$  and then  $\pi_2(p_2 | x_{1:m}, p_1)$  which are univariate truncated Gaussian distribution which can be sampled using the inverse cdf method. The ratio  $\pi(p_1, p_2 | x_{1:m}) / \pi_2(p_1, p_2 | x_{1:m})$  is upper bounded.

- (Geweke, 1989) proposed using the normal approximation to the binomial distribution

$$\pi_2(p_1, p_2 | x_{1:m}) \propto \exp\left(- (m_{1,1} + m_{1,2}) (p_1 - \hat{p}_1)^2 / (2\hat{p}_1 (1 - \hat{p}_1))\right) \\ \times \exp\left(- (m_{2,1} + m_{2,2}) (p_2 - \hat{p}_2)^2 / (2\hat{p}_2 (1 - \hat{p}_2))\right)$$

where  $\hat{p}_1 = m_{1,1} / (m_{1,1} + m_{1,2})$ ,  $\hat{p}_2 = m_{2,2} / (m_{2,2} + m_{2,1})$ . Then to simulate from this distribution, we simulate first  $\pi_2(p_1 | x_{1:m})$  and then  $\pi_2(p_2 | x_{1:m}, p_1)$  which are univariate truncated Gaussian distribution which can be sampled using the inverse cdf method. The ratio  $\pi(p_1, p_2 | x_{1:m}) / \pi_2(p_1, p_2 | x_{1:m})$  is upper bounded.

- A final one consists of using

$$\pi_3(p_1, p_2 | x_{1:m}) = \mathcal{B}e(p_1; m_{1,1} + 1, m_{1,2} + 1) \pi_3(p_2 | x_{1:m}, p_1)$$

where  $\pi(p_2 | x_{1:m}, p_1) \propto (1 - p_2)^{m_{2,1}} p_2^{m_{2,2}} \mathbb{I}_{p_2 \leq 1 - p_1}$  is badly approximated through  $\pi_3(p_2 | x_{1:m}, p_1) = \frac{2}{(1 - p_1)^2} p_2 \mathbb{I}_{p_2 \leq 1 - p_1}$ . It is straightforward to check that  $\pi(p_1, p_2 | x_{1:m}) / \pi_3(p_1, p_2 | x_{1:m}) \propto (1 - p_2)^{m_{2,1}} p_2^{m_{2,2}} / \frac{2}{(1 - p_1)^2} p_2 < \infty$ .

- Performance for  $N = 10,000$

Distribution	$\varphi_1$	$\varphi_2$	$\varphi_3$	$\varphi_4$	$\varphi_5$
$\pi_1$	0.748	0.139	3.184	0.163	2.957
$\pi_2$	0.689	0.210	2.319	0.283	2.211
$\pi_3$	0.697	0.189	2.379	0.241	2.358
$\pi$	0.697	0.189	2.373	0.240	2.358

- Performance for  $N = 10,000$

Distribution	$\varphi_1$	$\varphi_2$	$\varphi_3$	$\varphi_4$	$\varphi_5$
$\pi_1$	0.748	0.139	3.184	0.163	2.957
$\pi_2$	0.689	0.210	2.319	0.283	2.211
$\pi_3$	0.697	0.189	2.379	0.241	2.358
$\pi$	0.697	0.189	2.373	0.240	2.358

- Sampling from  $\pi$  using rejection sampling works well but is computationally expensive.  $\pi_3$  is computationally much cheaper whereas  $\pi_1$  does extremely poorly as expected.



# Optimal Normalized Importance Sampling

- For a given test function, one can minimize the normalized IS asymptotic variance using

$$q^{\text{opt}}(x) = \frac{|\varphi(x) - \mathbb{E}_{\pi}(\varphi)| \pi(x)}{\int_{\mathcal{X}} |\varphi(x) - \mathbb{E}_{\pi}(\varphi)| \pi(x) dx}$$

# Optimal Normalized Importance Sampling

- For a given test function, one can minimize the normalized IS asymptotic variance using

$$q^{\text{opt}}(x) = \frac{|\varphi(x) - \mathbb{E}_{\pi}(\varphi)| \pi(x)}{\int_{\mathcal{X}} |\varphi(x) - \mathbb{E}_{\pi}(\varphi)| \pi(x) dx}$$

- *Proof.*

$$\begin{aligned} & \int q(x) \frac{\pi^2(x)}{q^2(x)} (\varphi(x) - \mathbb{E}_{\pi}(\varphi))^2 dx \\ & \geq \left( \int q(x) \frac{\pi(x) |\varphi(x) - \mathbb{E}_{\pi}(\varphi)|}{q(x)} dx \right)^2 \\ & = \left( \int \pi(x) |\varphi(x) - \mathbb{E}_{\pi}(\varphi)| dx \right)^2 \end{aligned}$$

and this lower bound is attained for  $q^{\text{opt}}(x)$ .

# Optimal Normalized Importance Sampling

- For a given test function, one can minimize the normalized IS asymptotic variance using

$$q^{\text{opt}}(x) = \frac{|\varphi(x) - \mathbb{E}_{\pi}(\varphi)| \pi(x)}{\int_{\mathcal{X}} |\varphi(x) - \mathbb{E}_{\pi}(\varphi)| \pi(x) dx}$$

- *Proof.*

$$\begin{aligned} & \int q(x) \frac{\pi^2(x)}{q^2(x)} (\varphi(x) - \mathbb{E}_{\pi}(\varphi))^2 dx \\ & \geq \left( \int q(x) \frac{\pi(x) |\varphi(x) - \mathbb{E}_{\pi}(\varphi)|}{q(x)} dx \right)^2 \\ & = \left( \int \pi(x) |\varphi(x) - \mathbb{E}_{\pi}(\varphi)| dx \right)^2 \end{aligned}$$

and this lower bound is attained for  $q^{\text{opt}}(x)$ .

- This result is practically useless because it requires knowing  $\mathbb{E}_{\pi}(\varphi)$  but it suggests approximations.

# Effective Sample Size

- In statistics, we are usually not interested in a specific  $\varphi$  but in several functions and we prefer having  $q(x)$  as close as possible to  $\pi(x)$ .

# Effective Sample Size

- In statistics, we are usually not interested in a specific  $\varphi$  but in several functions and we prefer having  $q(x)$  as close as possible to  $\pi(x)$ .
- For flat functions, one can approximate the variance by

$$\text{var}(\mathbb{E}_{\hat{\pi}_N}(\varphi(X))) \approx (1 + \text{var}_q(w(X))) \frac{\text{var}(\mathbb{E}_{\pi}(\varphi(X)))}{N}.$$

# Effective Sample Size

- In statistics, we are usually not interested in a specific  $\varphi$  but in several functions and we prefer having  $q(x)$  as close as possible to  $\pi(x)$ .
- For flat functions, one can approximate the variance by

$$\text{var}(\mathbb{E}_{\hat{\pi}_N}(\varphi(X))) \approx (1 + \text{var}_q(w(X))) \frac{\text{var}(\mathbb{E}_{\pi}(\varphi(X)))}{N}.$$

- **Simple interpretation:** The  $N$  weighted samples are approximately equivalent to  $M$  unweighted samples from  $\pi$  where

$$M = \frac{N}{1 + \text{var}_q(w(X))} \leq N.$$

# Computing Ratio of Normalizing Constant

- However, we are often interested in estimating the ratio of normalizing constants

$$\frac{\int \pi^*(x) dx}{\int q^*(x) dx} = \int w^*(x) q(x) dx = \mathbb{E}_q [w^*(X)].$$

using

$$\mathbb{E}_{\hat{q}_N} [w^*(X)] = \frac{1}{N} \sum_{i=1}^N w^*(X^{(i)})$$

# Computing Ratio of Normalizing Constant

- However, we are often interested in estimating the ratio of normalizing constants

$$\frac{\int \pi^*(x) dx}{\int q^*(x) dx} = \int w^*(x) q(x) dx = \mathbb{E}_q [w^*(X)].$$

using

$$\mathbb{E}_{\hat{q}_N} [w^*(X)] = \frac{1}{N} \sum_{i=1}^N w^*(X^{(i)})$$

- It is unbiased and has variance

$$\text{var} [\mathbb{E}_{\hat{q}_N} [w^*(X)]] = \frac{\text{var}_q (w^*(X))}{N}.$$



- Clearly if you have  $q(x) = \pi(x)$  then

$$\text{var} [\mathbb{E}_{\hat{q}_N} [w^*(X)]] = 0$$

- Clearly if you have  $q(x) = \pi(x)$  then

$$\text{var} [\mathbb{E}_{\hat{q}_N} [w^*(X)]] = 0$$

- However if  $q(x) \neq \pi(x)$  then the estimate is simply

$$\mathbb{E}_{\hat{q}_N} [w^*(X)] = \frac{\int \pi^*(x) dx}{\int q^*(x) dx}.$$

- Clearly if you have  $q(x) = \pi(x)$  then

$$\text{var} [\mathbb{E}_{\hat{q}_N} [w^*(X)]] = 0$$

- However if  $q(x) = \pi(x)$  then the estimate is simply

$$\mathbb{E}_{\hat{q}_N} [w^*(X)] = \frac{\int \pi^*(x) dx}{\int q^*(x) dx}.$$

- **Open Question:** How could you come up with a good estimate of  $\int \pi^*(x) dx$  based on samples of  $\pi$ .

- IS is more powerful than you think.

- IS is more powerful than you think.
- Assume you have say to compute the importance weight

$$w(\theta) \propto \int f(x, z | \theta) dz$$

i.e. the likelihood is very complex and might not admit a closed-form expression.

- IS is more powerful than you think.
- Assume you have say to compute the importance weight

$$w(\theta) \propto \int f(x, z | \theta) dz$$

i.e. the likelihood is very complex and might not admit a closed-form expression.

- You do NOT need to compute  $w(\theta^{(i)})$  exactly, an unbiased estimate of it is sufficient.

# Limitations of Importance Sampling

- Consider the case where  $\mathcal{X} = \mathbb{R}^n$

$$\pi(\theta) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{\sum_{i=1}^n \theta_i^2}{2}\right)$$

and

$$q_\sigma(\theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{\sum_{i=1}^n \theta_i^2}{2\sigma^2}\right)$$

# Limitations of Importance Sampling

- Consider the case where  $\mathcal{X} = \mathbb{R}^n$

$$\pi(\theta) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{\sum_{i=1}^n \theta_i^2}{2}\right)$$

and

$$q_\sigma(\theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{\sum_{i=1}^n \theta_i^2}{2\sigma^2}\right)$$

- We have for any  $\sigma > 1$

$$w_\sigma(\theta) = \frac{\pi(\theta)}{q_\sigma(\theta)} = \sigma^n \exp\left(-\sum_{i=1}^n \frac{\theta_i^2}{2} \left(1 - \frac{1}{\sigma^2}\right)\right) \leq \sigma^n \text{ for any } \theta$$

and

$$\text{var}_{q_\sigma}\left(\frac{\pi(\theta)}{q_\sigma(\theta)}\right) = \sigma^n \sigma'^n - 1 \text{ with } \sigma'^2 = \frac{\sigma^2}{\sigma^2 - 1/2} > 1$$



# Limitations of Importance Sampling

- Consider the case where  $\mathcal{X} = \mathbb{R}^n$

$$\pi(\theta) = \frac{1}{(2\pi)^{n/2}} \exp\left(-\frac{\sum_{i=1}^n \theta_i^2}{2}\right)$$

and

$$q_\sigma(\theta) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{\sum_{i=1}^n \theta_i^2}{2\sigma^2}\right)$$

- We have for any  $\sigma > 1$

$$w_\sigma(\theta) = \frac{\pi(\theta)}{q_\sigma(\theta)} = \sigma^n \exp\left(-\sum_{i=1}^n \frac{\theta_i^2}{2} \left(1 - \frac{1}{\sigma^2}\right)\right) \leq \sigma^n \text{ for any } \theta$$

and

$$\text{var}_{q_\sigma}\left(\frac{\pi(\theta)}{q_\sigma(\theta)}\right) = \sigma^n \sigma'^n - 1 \text{ with } \sigma'^2 = \frac{\sigma^2}{\sigma^2 - 1/2} > 1$$

- Despite having a very good proposal then the variance of the weights increases exponentially fast with the dimension of the problem.

# Normalized Importance Sampling versus Rejection Sampling

- Given  $N$  samples from  $q$ , we estimate  $\mathbb{E}_\pi(\varphi(X))$  through IS

$$\mathbb{E}_{\hat{\pi}_N}^{\text{IS}}(\varphi(X)) = \frac{\sum_{i=1}^N w^*(X^{(i)}) \varphi(X^{(i)})}{\sum_{i=1}^N w^*(X^{(i)})}$$

or we “filter” the samples through rejection and propose instead

$$\mathbb{E}_{\hat{\pi}_N}^{\text{RS}}(\varphi(X)) = \frac{1}{K} \sum_{k=1}^K \varphi(X^{(i_k)})$$

where  $K \leq N$  is a random variable corresponding to the number of samples accepted.

# Normalized Importance Sampling versus Rejection Sampling

- Given  $N$  samples from  $q$ , we estimate  $\mathbb{E}_\pi(\varphi(X))$  through IS

$$\mathbb{E}_{\hat{\pi}_N}^{\text{IS}}(\varphi(X)) = \frac{\sum_{i=1}^N w^*(X^{(i)}) \varphi(X^{(i)})}{\sum_{i=1}^N w^*(X^{(i)})}$$

or we “filter” the samples through rejection and propose instead

$$\mathbb{E}_{\hat{\pi}_N}^{\text{RS}}(\varphi(X)) = \frac{1}{K} \sum_{k=1}^K \varphi(X^{(i_k)})$$

where  $K \leq N$  is a random variable corresponding to the number of samples accepted.

- We want to know which strategy performs the best.

- Define the artificial target  $\bar{\pi}(x, y)$  on  $\mathcal{X} \times [0, 1]$  as

$$\bar{\pi}(x, y) = \begin{cases} \frac{Cq^*(x)}{\int \pi^*(x) dx}, & \text{for } \left\{ (x, y) : x \in \mathcal{X} \text{ and } y \in \left[ 0, \frac{\pi^*(x)}{Cq^*(x)} \right] \right\} \\ 0 & \text{otherwise} \end{cases}$$

then

$$\int \bar{\pi}(x, y) dy = \int_0^{\frac{\pi^*(x)}{Cq^*(x)}} \frac{Cq^*(x)}{\int \pi^*(x) dx} dy = \pi(x).$$

- Define the artificial target  $\bar{\pi}(x, y)$  on  $\mathcal{X} \times [0, 1]$  as

$$\bar{\pi}(x, y) = \begin{cases} \frac{Cq^*(x)}{\int \pi^*(x) dx}, & \text{for } \left\{ (x, y) : x \in \mathcal{X} \text{ and } y \in \left[ 0, \frac{\pi^*(x)}{Cq^*(x)} \right] \right\} \\ 0 & \text{otherwise} \end{cases}$$

then

$$\int \bar{\pi}(x, y) dy = \int_0^{\frac{\pi^*(x)}{Cq^*(x)}} \frac{Cq^*(x)}{\int \pi^*(x) dx} dy = \pi(x).$$

- Now let us consider the proposal distribution

$$q(x, y) = q(x) \mathcal{U}_{[0,1]}(y) \text{ for } (x, y) \in \mathcal{X} \times [0, 1].$$

- Then rejection sampling is nothing but IS on  $\mathcal{X} \times [0, 1]$  where

$$w(x, y) = \frac{\bar{\pi}(x, y)}{q(x)\mathcal{U}_{[0,1]}(y)} = \begin{cases} \frac{C \int q^*(x) dx}{\int \pi^*(x) dx} & \text{for } y \in \left[0, \frac{\pi^*(x)}{Cq^*(x)}\right] \\ 0, & \text{otherwise.} \end{cases}$$

- Then rejection sampling is nothing but IS on  $\mathcal{X} \times [0, 1]$  where

$$w(x, y) = \frac{\bar{\pi}(x, y)}{q(x) \mathcal{U}_{[0,1]}(y)} = \begin{cases} \frac{C \int q^*(x) dx}{\int \pi^*(x) dx} & \text{for } y \in \left[0, \frac{\pi^*(x)}{Cq^*(x)}\right] \\ 0, & \text{otherwise.} \end{cases}$$

- We have

$$\mathbb{E}_{\hat{\pi}_N^{\text{RS}}}(\varphi(X)) = \frac{1}{K} \sum_{k=1}^K \varphi(X^{(i_k)}) = \frac{\sum_{i=1}^N w(X^{(i)}, Y^{(i)}) \varphi(X^{(i)})}{\sum_{i=1}^N w(X^{(i)}, Y^{(i)})}.$$

- Then rejection sampling is nothing but IS on  $\mathcal{X} \times [0, 1]$  where

$$w(x, y) = \frac{\bar{\pi}(x, y)}{q(x) \mathcal{U}_{[0,1]}(y)} = \begin{cases} \frac{C \int q^*(x) dx}{\int \pi^*(x) dx} & \text{for } y \in \left[0, \frac{\pi^*(x)}{Cq^*(x)}\right] \\ 0, & \text{otherwise.} \end{cases}$$

- We have

$$\mathbb{E}_{\hat{\pi}_N^{\text{RS}}}(\varphi(X)) = \frac{1}{K} \sum_{k=1}^K \varphi(X^{(i_k)}) = \frac{\sum_{i=1}^N w(X^{(i)}, Y^{(i)}) \varphi(X^{(i)})}{\sum_{i=1}^N w(X^{(i)}, Y^{(i)})}.$$

- Compared to standard IS, RS performs IS on an enlarged space.



- The variance of the importance weights from RS is higher than for standard IS:

$$\text{var}_q [w(X, Y)] \geq \text{var}_q [w(X)].$$

More precisely, we have

$$\begin{aligned} \text{var} [w(X, Y)] &= \text{var} [\mathbb{E} [w(X, Y) | X]] + \mathbb{E} [\text{var} [w(X, Y) | X]] \\ &= \text{var} [w(X)] + \mathbb{E} [\text{var} [w(X, Y) | X]]. \end{aligned}$$

- The variance of the importance weights from RS is higher than for standard IS:

$$\text{var}_q [w(X, Y)] \geq \text{var}_q [w(X)].$$

More precisely, we have

$$\begin{aligned} \text{var} [w(X, Y)] &= \text{var} [\mathbb{E} [w(X, Y) | X]] + \mathbb{E} [\text{var} [w(X, Y) | X]] \\ &= \text{var} [w(X)] + \mathbb{E} [\text{var} [w(X, Y) | X]]. \end{aligned}$$

- To compute integrals, Rejection sampling is inefficient and you should simply use IS.

- Like Rejection, IS is useful for small non-standard distributions but collapses for most “interesting” problems.

- Like Rejection, IS is useful for small non-standard distributions but collapses for most “interesting” problems.
- In both cases, the problem is to be able to design “clever” proposal distributions.

- Like Rejection, IS is useful for small non-standard distributions but collapses for most “interesting” problems.
- In both cases, the problem is to be able to design “clever” proposal distributions.
- Towards the end of this course, we will present advanced dynamic methods to address this problem.