

## Importance Sampling

Suppose your goal is to compute

$$E[f(x)] = \int_{-\infty}^{\infty} f(x) p(x) dx$$

Chebychev inequality (see Review) :

$$E[f(x)] \approx \frac{1}{N_c} \sum_{j=1}^{N_c} f(x_j), \quad x_j \sim p(x)$$

This works if  $N_c \rightarrow \infty$  limit, i.e.

$$\frac{1}{N_c} \sum f(x_j) \rightarrow E[f(x)] \text{ as } N_c \rightarrow \infty$$

In practice:  $N_c$  might be very large.

Example:

$$f(x) = \begin{cases} 0 & x < 4 \\ 1 & x \geq 4 \end{cases} \quad x \sim N(0, 1)$$

$$E[f(x)] = \int_{-\infty}^{\infty} f(x) p(x) dx = \int_{-\infty}^{\infty} f(x) \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x^2) dx$$

$$= \int_4^{\infty} \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x^2) dx \approx 3.17 \cdot 10^{-5}$$

Monte Carlo:

$$E[f(x)] \approx \frac{1}{N_c} \sum_{j=1}^{N_c} f(x_j) \quad x_j \sim N(0, 1)$$

Since "most"  $x_j$  are in  $(-2, 2)$

(see, e.g., Chebychev)

$f(x_j) = 0$  "most of the time".

↳ Unless  $N_c \rightarrow \infty$ , we get

$$E[f(x)] = \frac{1}{N_c} \sum f(x_j) = \frac{1}{N} \sum 0 = 0.$$

Why? Samples  $x_j \sim p(x)$  fall into regions where  $f(x)$  is small (zero) (-2,2)

Idea: focus samples onto region where  $f(x) \rightarrow \text{large}$ .

Define:

$$E_p[f(x)] = \int f(x) p(x) dx$$

↑ which distribution  
we average over

Then:

$$\begin{aligned} E_p[f(x)] &= \int f(x) p(x) dx \\ &= \int f(x) \frac{p(x)}{g(x)} g(x) dx \end{aligned}$$

assuming that  $g(x) = 0$  only when  $p(x) = 0$ ,  
i.e. support of  $g$  includes support of  $p$ .

$$= E_g \left[ f(x) \frac{p(x)}{g(x)} \right]$$

Mark Carlo: (Importance Sampling)

$$E_p[f(x)] = E_q\left[f(x) \frac{p(x)}{q(x)}\right] = \frac{1}{N_e} \sum_{j=1}^{N_e} f(x_j) \frac{p(x_j)}{q(x_j)}$$

$x_j \sim q(x_j)$

Terminology:  $p(x)$  → target distribution

$q(x)$  → proposal distribution

$$\omega(x) = \frac{p(x)}{q(x)} \rightarrow \text{weight function}$$

Algorithm: Draw  $N_e$  samples from proposal  $q$

• Compute

$$E[f(x)] \approx \frac{1}{N_e} \sum_{j=1}^{N_e} f(x_j) \omega(x_j)$$

This "biases" samples to fall into regions where  $f$  is large (important).

Back to example:

$$f(x) = \begin{cases} 0 & \text{if } x < 4 \\ 1 & \text{if } x \geq 4 \end{cases} \quad p(x) = N(0, 1)$$
$$q(x) = N(2, 1)$$

↳ This helps and gives better estimates for "small  $N_e$ ".

↳ You explore this in HW.

Question: How to pick a "good" proposal such that

$$\frac{1}{N_e} \sum_{j=1}^{N_e} f(x_j) \frac{P(x_j)}{g(x_j)} \quad x_j \sim g(x)$$

$\beta \sim$  "good" approximation of  $E_p[f(x)]$

for a "reasonable"  $N_e$ ?

Note: What "good" and "reasonable" means depends on the problem.

Recall the Review:

$$E_p[g(x)] = \frac{1}{N_e} \sum_{j=1}^{N_e} g(x_j) \quad x_j \sim p$$

↳ Error in estimate  $\Rightarrow$

$$\frac{\sigma(g(x))}{\sqrt{N_e}} \quad (\text{see Review})$$

$\sigma(g(x))$  is std. dev.  
of r.v.  $g(x)$ .

From Chebychev we get that

$$\gamma = \frac{1}{N_e} \sum \gamma_i; \quad \gamma_i \text{ i.i.d.}$$

$$P(|\gamma - E\gamma| > \kappa \frac{\sigma}{\sqrt{N_e}}) \leq \frac{1}{\kappa^2}$$

Applying this to importance sampling estimate:

$$\frac{\sigma(f(x) \frac{P(x)}{g(x)})}{\sqrt{N_e}}$$

is "error" in importance sampling -

$\Rightarrow$  we want that std. dev. of

$f(x) \frac{P(x)}{q(x)}$  is "small".

$\Rightarrow$  smallest it can be is zero, what happens when

$$f(x) \frac{P(x)}{q(x)} = \text{const.}$$

$\Rightarrow$  a good proposal distribution is such that

$$q(x) \propto P(x) f(x).$$

"Problem":  $q(x) \geq 0$ .

Definition: The optimal proposal distribution is

$$q^*(x) = \frac{|f(x)| p(x)}{E_p[|f(x)|]}$$

and minimizes  $\sigma(f(x) \frac{P(x)}{q(x)})$ .

(1.) Check that  $q^*$  is a distribution.

$$q^*(x) \geq 0 \quad \checkmark$$

$$\int_{-\infty}^{\infty} q^*(x) dx = \int \frac{|f(x)| p(x)}{E_p[|f(x)|]} dx = \frac{1}{E_p[|f(x)|]} E_p[|f(x)|] = 1 \quad \checkmark$$

$$(2) \sigma(f(x) \frac{P(x)}{q^*(x)}) \leq \sigma(f(x) \frac{P(x)}{q(x)})$$

Proof:

$$\text{Recall: } \text{Var}(x) = E(x^2) - E(x)^2$$

$$\sigma^2(\int p_{q^*}) = \text{var}(\int p_{q^*})$$

$$= E_{q^*}[(\int p_{q^*})^2] - \underbrace{E_{q^*}[\int p_{q^*}]^2}_{:= \mu = E[\int p_{q^*}]}$$

$$\Rightarrow \underbrace{\sigma_{q^*}^2(\int p_{q^*}) + \mu^2}_{=} = E_{q^*}[(\int p_{q^*})^2]$$

$$= \int \int(x)^2 \frac{p(x)^2}{g^*(x)^2} \sqrt{g^*(x)} dx = \int \frac{\int(x)^2 p(x)^2}{|\int(x)|^2 p(x)} dx \underbrace{E_p[|\int(x)|]}_{\int \frac{|\int(x)|^2}{|\int(x)|} p(x) dx = E_p[|\int(x)|]}$$

$$= E_p[|\int(x)|]^2 = E_g \left[ \left| \int(x) \right| \frac{p(x)}{g(x)} \right]^2$$

$$\begin{aligned} \text{var}(x) &= E(x^2) - E(x)^2 \text{ Some other} \\ &\Rightarrow E(x)^2 = E(x^2) - \text{var}(x) \text{ proposal} \end{aligned}$$

$$\Rightarrow E_g \left[ \left| \int(x) \right|^2 \frac{p(x)}{g(x)} \right] - \text{var} \left( \left| \int(x) \right| \frac{p(x)}{g(x)} \right)$$

$$\leq E_g \left[ \underbrace{\left| \int(x) \right|^2}_{= \int(x)^2} \frac{p(x)^2}{g(x)^2} \right] = \underbrace{\sigma_g^2(|\int|) + \mu^2}_{E(x^2) = \sigma^2 + E(x)^2}$$

$$\sigma_{q^*}^2(\int p_{q^*}) \leq \sigma_g^2(\int p_{q^*}) \quad \checkmark$$

## Problems:

(i) We don't know how to draw samples from  $p(x)$ .  
How do we draw samples from  $q^*(x) \propto p(x) |f(x)|$ ?

(ii) We need to know or compute  $\bar{E}_p[|f(x)|]$ .  
This seems difficult.

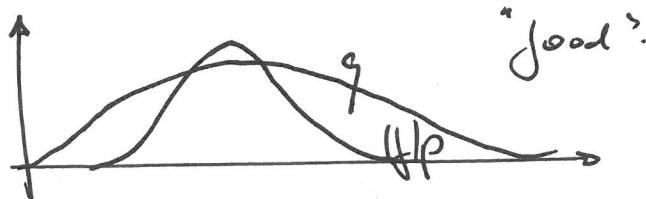
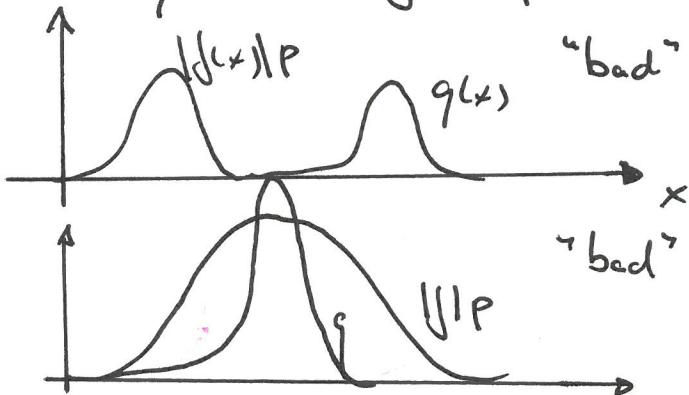
Lesson: - optimal proposal  $q^* = \frac{|f(x)| p(x)}{\bar{E}_p[|f(x)|]}$  usually  
hard to compute w.r.t.

- it is a useful tool to figure out what
- "good" proposal may look like.

(a)  $q(x)$  should be large where  $|f(x)| p(x)$  is large.

(b) avoid small  $q$  where  $|f(x)| p(x)$  is large.

In Cartoons:



## Another difficulty of importance sampling

$$E_p[f(x)] = E_g\left[f(x) \frac{P(x)}{g(x)}\right] \approx \frac{1}{N_c} \sum_{j=1}^{N_c} f(x_j) \frac{P(x_j)}{g(x_j)} \quad x_j \sim g(x).$$

weights :  $\omega(x_j) = \frac{P(x_j)}{g(x_j)}$

→ we need to know  $P(x)$ .

→ we often don't.

In PA :  $P_0(x) \rightarrow$  prior

$$y = Hx + \eta, \eta \sim N(0, R)$$

$$P(x|y) = \frac{P_0(x) P_e(y|x)}{P_y(y)}$$

We know  $P(x|y)$  up to a multiplicative constant, but t.c. const. is hard to compute.

$$P_y(y) = \int P_0(x) P_e(y|x) dx$$

another integral!  $(\int P(x|y) dx = 1)$

Help: Self-normalized importance sampling.

$$\omega(x_j) \propto \frac{P(x_j)}{g(x_j)} \quad \underbrace{\text{Note: Const. cannot}}_{j=1, \dots, N_c} \quad \text{depend on } x_j!$$

Define:  $\hat{\omega}(x_j) = \omega(x_j) / \sum_{j=1}^{N_c} \omega(x_j)$

IS independent of proportionality constant.

$$\omega(x_j) = C \frac{\hat{p}(x_j)}{\hat{q}(x_j)} \quad \text{where} \quad \hat{p}(x_j) \propto p(x_j)$$

$$\hat{q}(x_j) \propto q(x_j)$$

$$\hat{\omega}(x_j) = \frac{C \frac{\hat{p}(x_j)}{\hat{q}(x_j)}}{\sum_i C \frac{\hat{p}(x_i)}{\hat{q}(x_i)}} = \frac{\hat{p}(x_j)/\hat{q}(x_j)}{\sum_{j=1}^n \hat{p}(x_j)/\hat{q}(x_j)}$$

→ For self-normalized importance sampling, we only need to know  $q(x)$  and  $p(x)$  up to a multiplicative const.

→ This makes a feasible algorithm.

$$E_p(f(x)) \approx \sum_{j=1}^{N_e} f(x_j) \hat{\omega}(x_j)$$

Yet another difficulty of importance sampling.

Suppose we are interested in PDF, and the r.v.  $x|y$ . We want to compute, at least,  $E[x|y]$  and  $\text{cov}(x|y)$

$$\hookrightarrow \mu = E[x] = \int x p(x|y) dx$$

$$\sigma^2 = E[(x-\mu)^2] = \int (x-\mu)^2 p(x|y) dx$$

→ What is a good "general purpose" proposal to compute several expected values of r.v.  $x|y$ ?

Spoiler alert: there are no good general purpose proposals.

• Some are "OK".

→ To find out which ones are OK, we need to find a way to assess how good a proposal is.

Consider:

$$E_p [f(x)] \approx \sum_{j=1}^{N_e} f(x_j) w_j$$

$w_j$  are "normalized" weights i.e.  $\sum_{j=1}^{N_e} w_j = 1$

→ weights do not depend on  $f$ .

→ we can assess how good  $f$  is by looking at the weight.

Worst Case:  $w_j = 0$  for  $j \neq k$   
 $w_k = 1$

→ one weight is one, all others are zero

→ one sample hogs all the probability, but there is no guarantee it's a good one.

→ estimate:

$$E_p [f(x)] = \sum_{j=1}^{N_e} f(x_j) w(x_j) = f(x_k)$$

→ terrible estimate based on only one sample!

Example:

$$\mu = E_p(x) = \sum x_j \omega_j = x_k$$

$$\text{Var}(x) = E_p((x - \mu)^2) = \sum (x_j - \mu)^2 \omega_j = (x_k - \mu)^2 = x_k - x_k = 0$$

but var. might not be zero.

How can we capture this?

Define "effective sample size".

$$N_{\text{eff}} = \frac{N_e}{g}$$

Idea: IS estimate with  $N_e$  samples are "as good" as estimates based on samples from target  $p(x)$  with  $N_{\text{eff}}$  samples

$$\mu_p = \sum_{j=1}^{N_{\text{eff}}} f(x_j), x_j \sim p$$

$$\mu_g = \sum_{j=1}^{N_e} f(x_j) \omega_j, x_j \sim g$$

$$\mu_p \approx \mu_g$$

We define:

$$g^0 = \frac{E(\omega^2)}{E(\omega)^2}$$

This is a heuristic quantity.

There are other ways to define th.g. We don't yet know what is most suitable way to describe all this.

$\mathcal{G}$  captures all the worst and best cases.

best case :  $g(x) \propto p(x)$

$$\hookrightarrow \omega \propto \frac{g(x)}{p(x)} = \text{const.}$$

Then  $\mathcal{G} = \frac{E(\omega^2)}{E(\omega)^2} = \frac{\text{Var}(y_k) + E(\omega)^2}{E(\omega)^2} = 1 + \frac{\text{Var}(\omega)}{E(\omega)^2}$

$$\text{Var}(x) = E(x^2) - E(x)^2$$

$$= 1 \quad \text{if} \quad \text{Var}(\omega) = 0 \quad \text{and} \quad \omega = \text{Const.}$$

$$\Rightarrow N_{\text{eff}} = \frac{N_e}{\mathcal{G}} = N_e. \quad \checkmark$$

Worst Case:

$$\omega_1 = 0, \omega_j^- = 0 \quad \text{for } j \neq 1$$

Then  $\mathcal{G} = \frac{E(\omega^2)}{E(\omega)^2} \approx \frac{\frac{1}{N_e} \sum \omega_j^2}{\left( \frac{1}{N_e} \sum \omega_j \right)^2} = N_e \sum \omega_j^2 = N_e$

$$N_{\text{eff}} = \frac{N_e}{\mathcal{G}} = 1 \quad \checkmark$$

## Resampling:

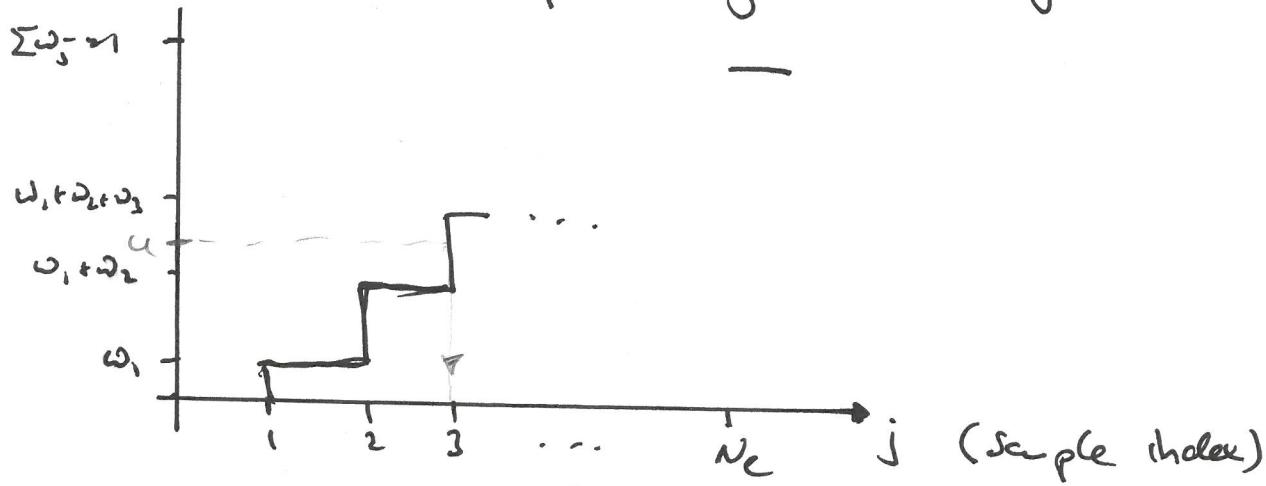
Idea: take a weighted ensemble  $\{x_j, \omega_j\}_{j=1}^{N_e}$  and transform it into an ensemble  $\{\tilde{x}_j\}_{j=1}^{N_e}$  with "equal weights" or "no weights".

$$\sum_{j=1}^{N_e} x_j \omega_j = \frac{1}{N_e} \sum_{j=1}^{N_e} x_j$$

One way of doing this is "resampling with replacement".

- Given a weighted ensemble  $\{x_j, \omega_j\}$

- Draw the PDF (probability distribution function)



- Draw uniform number between 0 and 1

If  $u < \omega_n$ : pick particle with index  $n$ .

- Do this  $N_e$  times.

You end up with an ensemble such that samples with large weights are repeated more often than samples with small weights.

And you have:

$$\begin{aligned} E_p[f(x)] &\approx \frac{1}{N_{\text{eff}}} \sum f(x_j) \quad x_j \sim p \\ &\approx \sum f(x_j) \tilde{\omega}_j \quad x_j \sim q \\ &\approx \frac{1}{N_e} \sum f(\tilde{x}_j), \quad \tilde{x}_j \sim q \end{aligned}$$

resampled ensemble.

$\Rightarrow$  You can use resampled ensemble in the same way as if you had an ensemble obtained by directly drawing from target distribution  $p$ .

$\Rightarrow$  Resampled ensemble can be used just like a "direct" ensemble for computing expected values

$\Rightarrow$  This is how importance sampling is most often used

$\Rightarrow$   $g$  and  $N_{\text{eff}}$  assess how good importance sampling works in this sense (i.e. independently of  $f(x)$ ).

### Histograms:

- Given an ensemble  $x_j \sim p(x_j) \quad j=1 \dots N_e$ ,  
e.g., obtained by importance sampling + resampling
- put samples  $x_j$  into "bins":

$$\exists x_j \in T_{jk} \Rightarrow x_j \text{ ends up in bin } k.$$

Draw height of bins (normalized to f.t. integral 1)



- ⇒ If you have a "good" ensemble, then histogram "looks like" probability distribution of  $x$ .
- ⇒ You can think of importance sampling as a technique to draw samples whose histogram looks like (converges, as  $N_e \rightarrow \infty$ ) the target distribution.