

Different Scenarios :

$$Z_p = \int \tilde{p}(x) dx$$

i.) Given x , compute $p(x)$, or $\tilde{p}(x)$, $p(x) = \frac{\tilde{p}(x)}{Z_p}$

ii) Sample from a distribution: Given $p(x)$, or more often $\tilde{p}(x)$,
generate samples $x_i \sim p(x)$, $i = 1, \dots, n$

- Evaluate statistics / moments: $\mathbb{E}_{x \sim p(x)} [f(x)] = \int f(x) p(x) dx$

⊛: Challenges arise when x is high-dimensional, $p(x)$ or $\tilde{p}(x)$ are complicated.

e.g. Bayesian inference: $\underbrace{p(\theta|D)}_{\hookrightarrow \text{draw } \theta \sim p(\theta|D)} \propto p(D|\theta)p(\theta)$

$$\underbrace{p(f_\theta(x^*) | x, y)}_{\text{predictive posterior}} = \int p(y|x, \theta) p(\theta|x, y) d\theta$$

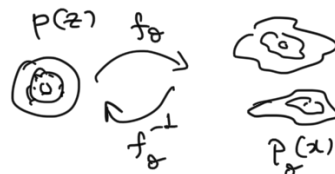
$$= \underbrace{\mathbb{E}_{\theta \sim p(\theta|D)} [p(y|x, \theta)]}_{\text{predictive posterior}}$$

Recall :

• Normalizing flow: given some samples x_i , estimate $p_\theta(x)$

$$\underbrace{p_\theta(x)}_{\text{Gaussian}} = \underbrace{p(z)}_{\text{Gaussian}} |\det \nabla_x f_\theta(x)| = p(f_\theta^{-1}(x)) |\det \nabla_x f_\theta^{-1}(x)|$$

$$z \sim \mathcal{N}(0, I), \quad x = f_\theta(z)$$



• Variational inference:

$$p(x) \approx q_{\phi}(x) \stackrel{\text{M.F.}}{=} \prod_{i=1}^n \mathcal{N}(x_i | \mu_i, \sigma_i^2)$$

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \operatorname{KL}[q_{\phi}(x) || p(x)]$$

Examples :

$$\text{i.) } \underset{p \in \mathcal{R}}{\mathbb{E}} [h(p)] \approx \frac{1}{S} \sum_{i=1}^S h(p_i) \xrightarrow{S \rightarrow +\infty}$$

ii) Making predictions using an ML model :

$$\begin{aligned} p(f(x^*) | \mathcal{D}) &= \mathbb{E}_{\theta \sim p(\theta | \mathcal{D})} [p(f(x^*) | \mathcal{D}, \theta)] \\ &\approx \frac{1}{n} \sum_{\theta_i \stackrel{\text{i.i.d.}}{\sim} p(\theta | \mathcal{D})} p(f(x^*) | \mathcal{D}, \theta_i), \end{aligned}$$

Why sampling?

i.) Allows us to compute expectations in high-dimensions $\left\{ \begin{array}{l} \text{estimate} \\ \text{statistical} \\ \text{posterior} \end{array} \right.$

$$p(X \in A) = \mathbb{E} [\mathbb{1}_{\{X \in A\}}], \quad \mathbb{1}_{\{X \in A\}} = \begin{cases} 1, & \text{if } X \in A \\ 0, & \text{if } X \notin A \end{cases}$$

ii) Allows us to compute intractable sums or integrals

Pros : 1.) Easy to implement / understand

2.) General purpose

3.) Well understood asymptotical theoretical guarantees.

... (but they may be inefficient in practice)

- Cons:
- 1.) They're too simple, ... often used inappropriately.
 - 2.) They tend to slow compared to deterministic approaches
 - 3.) It can be difficult to assess their performance.

Monte Carlo Approximation

Goal: Approximate an expectation using samples:

$$\mathbb{E}_{x \sim p(x)} [f(x)], \quad x \in \mathbb{R}^d.$$

Definition: If $x_1, x_2, \dots, x_n \stackrel{i.i.d.}{\sim} p(x)$, then:

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n f(x_i) : \text{a basic Monte Carlo (MC) estimator}$$

$$\approx \int f(x) p(x) dx$$

Remarks:

$$1.) \quad \mathbb{E}[\hat{\mu}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(x_i)] \xrightarrow{n \rightarrow +\infty} \mathbb{E}[f(x)]$$

hence $\hat{\mu}_n$ is an unbiased estimator

$$2.) \quad \hat{\mu}_n \xrightarrow[n \rightarrow +\infty]{P} \mathbb{E}[f(x)], \quad \forall \varepsilon > 0, \quad P(|\hat{\mu}_n - \mathbb{E}[f(x)]| < \varepsilon) \xrightarrow[n \rightarrow +\infty]{} 1$$

hence $\hat{\mu}_n$ is a consistent estimator
($x_i \stackrel{i.i.d.}{\sim} p(x)$)

$$3.) \quad \text{Var}[\hat{\mu}_n] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[f(x_i)] = \frac{1}{n} \underbrace{\text{Var}[f(x)]}_{\text{unknown}} \xrightarrow[n \rightarrow +\infty]{} 0$$

$$\text{std}[\hat{\mu}_n] \sim C \cdot \underbrace{O\left(\frac{1}{\sqrt{n}}\right)}_{\sim d!}$$

does not depend on n .

* Practical limitation:

MC approximation relies on the fact that we can efficiently

sample $x \sim p(x)$.

Importance Sampling (not a sampling method)

i.) Use I.S. to approximate $\mathbb{E}_{x \sim p(x)} [f(x)]$ when sampling from $p(x)$ is not tractable

ii.) Use I.S. to improve the accuracy/convergence rate of M.C. even if we could efficiently sample from $p(x)$

Setup: Assume that $p(x)$ is a density, i.e. $\int p(x) dx = 1$

$$\mathbb{E}_{x \sim p(x)} [f(x)] = \int f(x) p(x) dx \stackrel{\text{I.S.}}{=} \int \left[f(x) \frac{p(x)}{q(x)} \right] q(x) dx$$

$$= \mathbb{E}_{x \sim q(x)} \left[f(x) \frac{p(x)}{q(x)} \right] \approx \frac{1}{n} \sum_{i=1}^n f(x_i) w(x_i),$$

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

importance weight

$x_i \stackrel{\text{i.i.d.}}{\sim}$

$q(x)$

proposal distribution.

- i.) Easy to evaluate. { e.g. }
ii.) Easy to sample from. { e.g. }

Remarks:

$$\mathbb{E} \left[\hat{\mu}_n^{\text{IS}} \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{x \sim q(x)} \left[f(x_i) \frac{p(x_i)}{q(x_i)} \right] = \frac{1}{n} \sum_{i=1}^n \int f(x_i) \frac{p(x_i)}{q(x_i)} q(x_i) dx_i$$

$$x \sim q(x)$$

$$= \hat{\mu}_n^{\text{m.c.}} \text{ which is unbiased.}$$

$$\text{Var}[\hat{\mu}_n^{\text{IS}}] = \frac{1}{n} \text{Var}\left[f(x) \underbrace{\frac{p(x)}{q(x)}}_{C \cdot \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)}\right]$$

could be made smaller by appropriately choosing $q(x)$

In fact, it is straightforward to solve for the optimal $q^*(x)$:

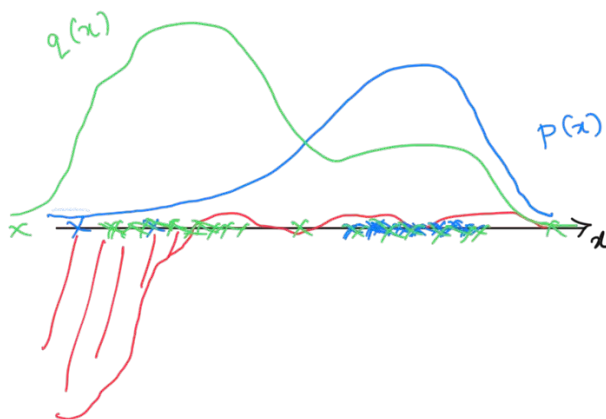
$$\boxed{q^*(x) = \min_q \text{Var}[\hat{\mu}_n^{\text{IS}}]}$$

How to choose a good importance sampling proposal distrib

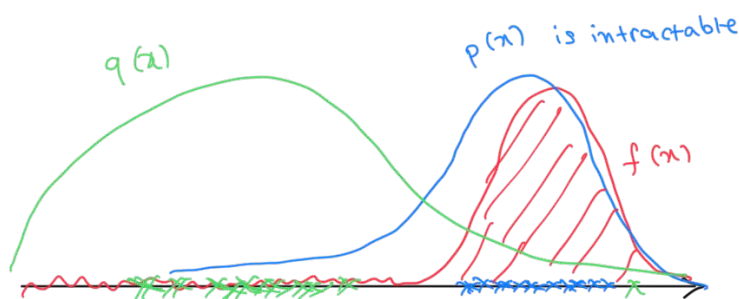
Example: $x \in \mathcal{X}$: investment
 $f(x)$: return

Approximate the expected return

$$\begin{aligned} \mathbb{E}_{x \sim p(x)}[f(x)] &= \int f(x) p(x) dx \\ &\approx \frac{1}{n} \sum_{i=1}^n f(x_i) \\ &\quad x_i \stackrel{\text{i.i.d.}}{\sim} p(x) \end{aligned}$$



$$\boxed{\approx \frac{1}{n} \sum_{i=1}^n f(x_i) w(x_i), \quad x_i \sim q(x)}$$



Sampling from a "wrong" $q(x)$ can make things go horribly wrong!

Choose $q(x)$ to be large when $|f(x)|p(x)$ is large.

Caution: Choosing $q(x)$ in high-dimensions is very difficult.

- It is difficult to assess how good or bad our estimator.

