

Sampling Methods :

Different cases :

1. Given $p(x)$, draw samples $x \sim p(x)$

e.g. $x \sim \mathcal{N}(\mu, \Sigma)$, $x = \mu + LZ$, $z \sim \mathcal{N}(0, I)$
 \mathbb{R}^d $d \times 1$ $d \times d$

$$\Sigma = LL^T \quad \begin{array}{l} \text{Cholesky} \\ \text{decomposition} \end{array}$$

\hookrightarrow Cholesky factor

e.g. $p(\theta|D) \propto p(D|\theta)p(\theta)$

we want to generate $\theta \sim p(\theta|D)$.

2. Given x_1, x_2, \dots, x_n , $x_i \in \mathbb{R}^d$, learn $p(x)$

3. Estimate statistics, e.g. given a r.v $x \sim p(x)$:

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

4. Perform Bayesian inference.

* Many sums or integrals can be written as expectations :

e.g. • marginal likelihood : $p(D) = \int p(D|\theta)p(\theta)d\theta =$
 $= \mathbb{E}_{\theta \sim p(\theta)}[p(D|\theta)]$

• predictive posterior distribution :

$$p(y^*|x^*, D) = \int p(y^*|x^*, D, \theta)p(\theta|D)d\theta$$
$$= \mathbb{E}_{\theta \sim p(\theta|D)}[p(y^*|x^*, D, \theta)]$$

Example :

Q1 : What is the avg height of students in ENM360 ?

$$\mu = 1.7 \dots 1 < 1.1$$

ans: $\mathbb{E}[L^u] := \frac{1}{|G|} \sum_{p \in G} L(p)$

Q2: What is the avg height of people in Center city?

ans: $\mathbb{E}[h] \approx \frac{1}{S} \sum_{i=1}^S h(p_i)$

Monte Carlo approximation:

Goal: Approximate an expectation/integral using samples.

$$\mathbb{E}[f(x)] = \int f(x) p(x) dx$$

$x \sim p(x)$ intractable

$x \in \mathbb{R}^d, d \geq 1$

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

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

(this is just the sample average)

Remarks:

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

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

2.) $\hat{\mu}_n \xrightarrow{P} \mathbb{E}[f(x)]$ as $n \rightarrow \infty$, convergence in probability

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

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

because x_i are i.i.d.

$$|\hat{\mu}_n - \mathbb{E}[f(x)]|^2 = \overset{0}{\text{bias}} + \text{var} \xrightarrow{n \rightarrow \infty} \frac{1}{n} \text{Var}[f(x)]$$

$$\xrightarrow{n \rightarrow \infty} \frac{1}{\sqrt{n}} \text{std}[f(x)]$$

Therefore, $\hat{\mu}_n$ converges to $\mathbb{E}[f(x)]$ at a rate $\mathcal{O}(\frac{1}{\sqrt{n}})$

Q: How many samples do we need in practice?

Note: Despite the fact that we know this rate of

1.) convergence, it may be very difficult what the actual error is in practice, because we don't

know the variance $\text{Var}[f(x)] = \int (f(x) - \mathbb{E}[f(x)])^2 p(x) dx$

2.) Practical limitation: One needs to be able to efficiently generate i.i.d. samples x_i from $p(x)$.