

Vanilla Monte Carlo

Vanilla Monte Carlo

- Consider evaluating the statistical expectation $\text{Exp}_p[f(x)]$

$$I = \text{Exp}_p[f(x)] = \int_x f(x) p(x) dx$$

Note that any general integral $\int f(x) dx$ can be written as a statistical expectation.

$$\int f(x) dx = \int \frac{f(x)}{p(x)} p(x) dx = \text{Exp}_p\left[\frac{f(x)}{p(x)}\right]$$

- We can estimate the integral using Riemann discretisation, taking evenly (deterministic) spaced pts in the domain of x (pts are separated by Δ)

$$\hat{I} = \sum_{i=1}^N \Delta_i f(x_i) p(x_i)$$

However, this has a few drawbacks

- It does not scale well w/ the dim. of X (req. many pts, hard to discretise / choose grid)
- The discretisation error is hard to characterise
- Wasteful in regions w/ small contributions to the integral.
- We can alternatively estimate the integral using a Vanilla Monte Carlo w/ N samples

$$\hat{I} = \frac{1}{N} \sum_{i=1}^N f(x_i)$$

where x_i are drawn from $p(x)$, i.e. $x \sim p$.

- * Note that the MC estimate \hat{I} is a RV \rightarrow has properties like expectation and variance.

Expectation and variance of the MC estimator

- The expectation of the MC estimator \hat{I} is given by

$$E[\hat{I}] = \frac{1}{N} \sum_{i=1}^N \text{Exp}_p[f(x_i)] = \frac{1}{N} \sum_{i=1}^N I = I$$

→ The MC estimator is unbiased ($E[\hat{I}] = I$)

- The variance of the MC estimator \hat{I} is given by

$$\text{Var}[\hat{I}] = \frac{1}{N^2} \sum_{i=1}^N \text{Var}[f(x_i)] = \frac{1}{N^2} \sum_{i=1}^N \sigma_f^2 = \frac{\sigma_f^2}{N}$$

→ The variance of the MC estimator decreases linearly w/ the no. of samples N but also depends on σ_f^2 so if it is exorbitant, even large values of N will not help.

- * since $E[\hat{I}] = I$ and $\lim_{N \rightarrow \infty} \text{Var}[\hat{I}] = 0$, the MC estimator is consistent

For Personal Use Only -bkwk2

Limit theorems for MC estimate

- The CLT states if $Y_1, Y_2 \dots$ be a sequence of iid RV w/ mean μ and variance $\sigma^2 < \infty$, and defining $S_n = \frac{1}{n} \sum_{i=1}^n Y_i$, then $\frac{S_n - \mu}{\sqrt{n}} \xrightarrow{D} N(0, 1)$ as $n \rightarrow \infty$

↳ W/o loss of generality, set $\mu = 0$. so $Z_n = \frac{S_n}{\sigma \sqrt{n}} = \frac{S_n - \mu}{\sigma \sqrt{n}}$

consider the mgf of Z_n , $M_{Z_n}(t) = E[\exp(-t Z_n)]$

$$\begin{aligned} M_{Z_n}(t) &= E[\exp(-t Z_n)] = E[(\exp(-t \frac{S_n}{\sigma \sqrt{n}}))^n] = (E[(1 - \frac{t}{\sigma \sqrt{n}} S_n + \frac{t^2}{\sigma^2 n} S_n^2 + \dots)])^n \\ &= (1 - \frac{t}{\sigma \sqrt{n}} E[S_n] + \frac{t^2}{\sigma^2 n} E[S_n^2] + \dots)^n = (1 + \frac{t^2}{2n} + \dots)^n \end{aligned}$$

Noting that $\lim_{n \rightarrow \infty} (1 + \frac{x}{n})^n = e^x$, $\lim_{n \rightarrow \infty} M_{Z_n}(t) = \lim_{n \rightarrow \infty} (1 + \frac{t^2}{2n} + \dots)^n = \exp(\frac{t^2}{2})$

which is the mgf of $N(0, 1)$.

- The SLN states if $Y_1, Y_2 \dots$ be a sequence of iid RV w/ mean μ , and defining $S_n = \frac{1}{n} \sum_{i=1}^n Y_i$, then $S_n \xrightarrow{P} \mu$ w/ probability 1 as $n \rightarrow \infty$.

- Set $Y_i = f(x_i)$, then the MC estimate \hat{I} becomes $\hat{I} = \frac{1}{N} \sum_{i=1}^N f(x_i) = \frac{1}{N} \sum_{i=1}^N Y_i = S_N$

Applying CLT, for large N , $\hat{I} \approx N(\mu, \frac{\sigma^2}{N})$

Applying SLN, f is guaranteed to converge to μ w/ probability 1 as $N \rightarrow \infty$

Variance reduction

Importance sampling (IS)

- To estimate the integral $E_{p(x)}[f(x)]$, we can use a proposal distribution $q(x)$

to sample from instead of $p(x)$

$$I = E_{p(x)}[f(x)] = \int_x f(x) p(x) dx = \int_x f(x) \frac{p(x)}{q(x)} q(x) dx = E_q[f(x) \frac{p(x)}{q(x)}]$$

- Given a proposal function $q(x)$ s.t. $q(x) > 0$, whenever $p(x) > 0$, the IS estimate is

$$\hat{I} = \frac{1}{N} \sum_{i=1}^N f(x_i) \frac{p(x_i)}{q(x_i)}$$

where x_i are drawn from $q(x)$, i.e. $x_i \sim q$

- The IS estimator is unbiased, $E[\hat{I}] = \frac{1}{N} \sum_{i=1}^N E_{q(x)}[f(x_i) \frac{p(x_i)}{q(x_i)}] = \frac{1}{N} \sum_{i=1}^N I = I$

- Define $F(x) = \frac{f(x)p(x)}{q(x)}$. The Vanilla MC estimate and IS estimate have errors:

$$\hookrightarrow \text{Vanilla MC : } \sigma_f^2 = \text{Var}[f(x)] = E_p[f(x)] - E_p[f(x)]^2 = E_p[f(x)^2] - I^2$$

$$\hookrightarrow \text{IS : } \sigma_I^2 = \text{Var}[f(x)] = E_q[f(x)] - E_q[f(x)]^2 = E_q[f(x)^2] - I^2$$

$$\Rightarrow \text{Diff in errors } \sigma_f^2 - \sigma_I^2 = E_p[f(x)^2] - E_q[f(x)^2] = \int f(x)^2 p(x) dx - \int \frac{f(x)^2 p(x)}{q(x)} p(x) dx$$

- By calculus of variations, the biggest decrease in error is achieved when $q(x) \propto f(x)p(x)$

$$q_{opt}(x) = \frac{f(x)p(x)}{\int f(x)p(x) dx}$$

intuitively, we want sample regions w/ large $f \rightarrow \uparrow$ contributions

$$\text{where we have } E_{q_{opt}}[(f(x) \frac{p(x)}{q(x)})^2] = \int_x \frac{f(x)^2 p(x)^2}{q(x)} dx = \left[\int_x f(x) p(x) dx \right]^2 = I^2 \rightarrow \sigma_I^2 = 0$$

- * IS is useful for reducing estimator variance and is useful when we cannot draw samples from $p(x)$. (we don't even req. normalized distributions of p, q)

For Personal Use Only -bkwk2

control variates

- say we want to estimate $E[f(x)]$ and we have additional info in the form of another function $g(x)$ where we know the value of $E[g(x)]$.
- Denoting the MC estimate of

$$\hat{E}_f = \frac{1}{N} \sum_i f(x_i)$$

$$\hat{E}_g = \frac{1}{N} \sum_i g(x_i)$$

An unbiased estimate follows as

$$\hat{I} = \hat{E}_f^c = \hat{E}_f + c(\hat{E}_g - E[g(x)])$$

where c controls the quality of the estimate.

- the estimator is unbiased $E[\hat{E}_f^c] = E[\hat{E}_f] + c(E[\hat{E}_g] - E[g(x)]) = E_f$
- the estimator has variance $\text{Var}[\hat{E}_f^c] = \text{Var}[\hat{E}_f] + c^2 \text{Var}[\hat{E}_g] + 2c\text{Cov}[\hat{E}_f, \hat{E}_g]$

- the estimator has min. variance when

$$\frac{d\text{Var}[\hat{E}_f^c]}{dc} = 2c\text{Var}[\hat{E}_g] + 2\text{Cov}[\hat{E}_f, \hat{E}_g] = 0 \rightarrow c_{opt} = -\frac{\text{Cov}[\hat{E}_f, \hat{E}_g]}{\text{Var}[\hat{E}_g]}$$

The variance of the estimator is thus

$$\text{Var}[\hat{E}_f^{opt}] = \text{Var}[\hat{E}_f] - \frac{\text{Cov}[\hat{E}_f, \hat{E}_g]^2}{\text{Var}[\hat{E}_g]}$$

i.e. if $\text{Cov}[\hat{E}_f, \hat{E}_g] \neq 0$, the variance of the estimator is reduced.

- several control variates can be used at the same time

$$\hat{I} = \hat{E}_f^c = \hat{E}_f + \sum_{i=1}^M c_i (\hat{E}_{g_i} - E[g_i(x)])$$

- the estimator is unbiased $E[\hat{E}_f^c] = E[\hat{E}_f] + \sum_{i=1}^M c_i (E[\hat{E}_{g_i}] - E[g_i(x)]) = E_f$

$$\text{the estimator has variance } \text{Var}[\hat{E}_f^c] = \text{Var}[\hat{E}_f] + 2\mathbf{C}^T \mathbf{b} + \mathbf{C}^T \mathbf{C}$$

where $\mathbf{C} \in \mathbb{R}^M$ has elements c_i , $\mathbf{b} \in \mathbb{R}^M$ has elements $\text{Cov}[\hat{E}_f, \hat{E}_{g_i}]$, $\mathbf{C}^T \mathbf{C} \in \mathbb{R}^{MM}$ has elements $\text{Cov}[\hat{E}_{g_i}, \hat{E}_{g_j}]$

- the estimator has min. variance when

$$\frac{d\text{Var}[\hat{E}_f^c]}{d\mathbf{C}} = 2\mathbf{b} + 2\mathbf{C}^T \mathbf{C} = 0 \rightarrow c_{opt} = -\mathbf{C}^T \mathbf{b}$$

The variance of the estimator is thus

$$\text{Var}[\hat{E}_f^{opt}] = \text{Var}[\hat{E}_f] - \mathbf{b}^T \mathbf{C}^T \mathbf{b}$$

MEASURE THEORYPower set and σ -algebras

- Let S be a set. The power set P on S is the collection of all subsets of S , i.e.

$$P = \{A : A \subseteq S\}$$

e.g. $S = \{a, b, c\} \rightarrow P = \{\emptyset, S, \{a\}, \{b\}, \{c\}, \{a, b\}, \{b, c\}, \{a, c\}\}$.

- Let S be a set. A σ -algebra F on S is a collection of subsets of S that satisfies

(i) $\emptyset, S \in F$

(ii) $A \in F \Rightarrow A' = S \setminus A \in F$

(iii) For any countable seq. (A_n) , $n \in \mathbb{N}$ in F , we have that $\bigcup_{n \in \mathbb{N}} A_n \in F$.

\Rightarrow Any set $A \in F$ is a measurable set, and

The pair (S, F) is a measurable space

- Let S be a set, and $M \subseteq P(S)$ be a collection of subsets of S . The σ -algebra

generated by M , (the smallest σ -algebra $F(S)$ on S that contains M), $\sigma(M)$ is

$$\sigma(M) = \{M \subseteq S : M \in F \text{ } \forall \text{ } \sigma\text{-algebra } F \text{ that contain } M\}.$$

e.g.: $S = \{a, b, c, d\}$, $M = \{\{a\}, \{b\}\} \rightarrow \sigma(M) = \{\emptyset, S, \{a\}, \{b\}, \{a, b\}, \{c, d\}, \{a, c, d\}, \{b, c, d\}\}$,

e.g.: $S = \mathbb{Z}$, $M = \{\{x\} : x \in \mathbb{Z}\} \rightarrow \sigma(M) = P(S)$.

- Let $S = \mathbb{R}$, $M = \{U \subseteq \mathbb{R} : U \text{ is open}\}$, then $\sigma(M)$ is the Borel σ -algebra $B(\mathbb{R})$, i.e. it is the smallest σ -algebra that contains all open sets.

Measure and measure space

- A measure on a measurable space (S, F) is a function $M : F \rightarrow [0, \infty]$ that satisfies

(i) $M(\emptyset) = 0$

(ii) For any disjoint seq. (A_n) in F , we have that $M\left(\bigcup_{n \in \mathbb{N}} A_n\right) = \sum_{n \in \mathbb{N}} M(A_n)$

\Rightarrow We call (S, F, M) a measure space. The measure space is finite if $M(S) < \infty$.

- Consider the Borel σ -algebra $B(\mathbb{R})$. We can define the Lebesgue measure $M_L(A)$, $\lambda(A)$ of set A

$$M_L(A) = \lambda(A) = \inf \left\{ \sum_{n \in \mathbb{N}} b_n - a_n : a_n < b_n \text{ are st. } A \subseteq \bigcup_{n \in \mathbb{N}} [a_n, b_n] \right\}$$

i.e. for the simple case $A = [a, b]$, $M_L(A) = b - a$.

- Consider the Borel σ -algebra $B(\mathbb{R})$. We can define the counting measure $M_C(A)$ of set A

$$M_C(A) = \begin{cases} |A \cap \mathbb{N}| & \text{if } A \cap \mathbb{N} \text{ is finite} \\ \infty & \text{otherwise} \end{cases}$$

i.e. the no. of elements in A . It is trivial to show that $(\mathbb{R}, B(\mathbb{R}), M_C)$ is a measure space

For Personal Use Only -bkwk2

Measurable functions

- A function $f: X \rightarrow Y$ between two measurable spaces (X, Σ_X) and (Y, Σ_Y) is (Σ_X, Σ_Y) -measurable if $\forall A \in \Sigma_Y, f^{-1}(A) \in \Sigma_X$, where $f^{-1}(A) = \{x \in X : f(x) \in A\}$.
(For general functions f , $f^{-1}(A)$ is not guaranteed to be in Σ_X).
⇒ If $X = Y = \mathbb{R}$, and $\Sigma_X = \Sigma_Y = \mathcal{B}(\mathbb{R})$, then such functions f are Borel measurable functions.
- Given a measure space (X, Σ_X, μ) , a measure space (Y, Σ_Y) and a measurable function $f: X \rightarrow Y$. The pushforward measure $f_*\mu: \Sigma_Y \rightarrow [0, \infty]$ is defined as
$$f_*\mu(A) = \mu(f^{-1}(A))$$
This means a measure space (Y, Σ_Y, μ) and measurable function $f: X \rightarrow Y$ together induce a measure on (Y, Σ_Y) .

Probability spaces

- We call a measure space $(\Omega, \mathcal{F}, \mu)$ a probability space if $\mu(\Omega) = 1$. We usually use the notation (Ω, \mathcal{F}, P) .
- A probability space (Ω, \mathcal{F}, P) represents an experiment.
 - ↪ Ω is the sample space. Any element we'll is an outcome of the experiment.
 - ↪ \mathcal{F} is a family of events (each event describes whether some phenomenon happened or not)
 - ↪ P is the measure assigning probabilities to each event.

(Intuitively, for a measurable set $A \in \mathcal{F}$, $P(A)$ represents the probability that a randomly selected outcome will belong to A .)
- For a Borel measurable function f , if $(\Omega, \mathcal{F}, \mu)$ is a probability space, f is called a random variable (RV), the pushforward measure $f_*\mu$ is the law of random variable f .

Almost everywhere and almost surely

- If turns out that whatever happens on a set of measure 0 simply does not matter.
- If a property of a measurable function holds everywhere, except on some set of measure, then we say the property holds almost everywhere (a.e.).
- When we work w/ probability spaces (Ω, \mathcal{F}, P) , we say the property holds almost surely (a.s.).

For Personal Use Only -bkwk2

Riemann and Lebesgue integration

Convergence of functions

- A sequence of functions f_1, f_2, \dots can converge to a function f either in a pointwise manner (each pt. could have diff. rates of convergence) or a uniform manner

- For pointwise convergence, we have $\lim_{n \rightarrow \infty} f_n(x) = f(x) \quad \forall x \quad [\lim_{n \rightarrow \infty} f_n = f]$

Formally, $\forall \epsilon > 0, \exists$ a natural no. $N = N(\epsilon, x)$ s.t. $n \geq N \Rightarrow |f(x) - f_n(x)| < \epsilon$

→ e.g. consider $f_n(x) = x^n, x \in [0, 1]$, $f_n(x)$ converges pointwise to $f(x)$

$$f(x) = \lim_{n \rightarrow \infty} f_n(x) = \begin{cases} 0 & x \in [0, 1) \\ 1 & x = 1 \end{cases}$$

Note $f_n(x)$ is continuous, $f(x)$ is not → pointwise convergence doesn't preserve continuity.

- For uniform convergence, we have $\sup_x |f(x) - f_n(x)| \rightarrow 0 \quad [f_n \xrightarrow{u} f]$

Formally, $\forall \epsilon > 0, \exists$ a natural no. $N = N(\epsilon)$ s.t. $n \geq N \Rightarrow |f(x) - f_n(x)| < \epsilon$

→ e.g. consider $f_n(x) = \frac{x}{n}, x \in [0, 1]$, $f_n(x)$ converges uniformly to $f(x) = 0$.

Riemann integration

- In Riemann integration, we approx. the integrand w/ step functions comprised of rectangles.

- We say that $s: [a, b] \rightarrow \mathbb{R}$ is a step function, if for some partition $a = t_0 < t_1 < \dots < t_N = b$

$$s(x) = \sum_{i=1}^N c_i \mathbb{I}_{[t_{i-1}, t_i)}(x)$$

where $c_i \in \mathbb{R}$ are constants.

- We define the integral of a step function as

$$\int_a^b s(x) dx = \sum_{i=1}^N c_i (t_i - t_{i-1})$$

- Let $f: [a, b] \rightarrow \mathbb{R}$ be a general function, define:

$$I_- = \sup_s \int_a^b s(x) dx \quad I_+ = \inf_s \int_a^b s(x) dx$$

If the bound from above / below are equal, i.e. $I_- = I_+$, then f is integrable and its integral is defined as $I = \int_a^b f(x) dx$

- There are a few fundamental problems w/ the Riemann integral:

→ It is difficult to integrate over infinite-dim. spaces — hard to define on interval.

→ The exchange of limit operations req. uniform convergence of f_n (which is stringent)

$$\lim_{n \rightarrow \infty} \int_a^b f_n(x) dx = \int_a^b \lim_{n \rightarrow \infty} f_n(x) dx$$

$$\text{proof: } \left| \lim_{n \rightarrow \infty} \int_a^b f_n(x) dx - \int_a^b f(x) dx \right| = \left| \lim_{n \rightarrow \infty} \int_a^b f_n(x) - f(x) dx \right|$$

$$\leq \lim_{n \rightarrow \infty} \int_a^b |f_n(x) - f(x)| dx \leq \lim_{n \rightarrow \infty} (b-a) \sup_x |f_n(x) - f(x)| = 0$$

* THIS IS NOT TRUE FOR POINTWISE CONVERGENCE

uniform convergence means

$$\lim_{n \rightarrow \infty} \sup_x |f_n(x) - f(x)| = 0.$$

For Personal Use Only -bkwk2

Lebesgue Integration

- In Lebesgue integration, rectangles of the same height can be joined together into one shape. The area covered by the shape is the height times the measure of the set of intervals.
- Let (Ω, \mathcal{F}) be a measurable space. $s: \Omega \rightarrow \mathbb{R}$ is a simple function if for some measurable sets $E_1, \dots, E_N \in \mathcal{F}$,

$$s(x) = \sum_{i=1}^N c_i I_{E_i}(x)$$

where $c_i \in \mathbb{R}$ are constants

- If μ is a measure on (Ω, \mathcal{F}) , the integral of a simple function wrt μ is defined as

$$\int_{\Omega} s(x) \mu(dx) = \sum_{i=1}^N c_i \mu(E_i)$$

- Let $(\Omega, \mathcal{F}, \mu)$ be a measure space and $f: \Omega \rightarrow [0, \infty]$ a non-negative measurable function. The integral of f wrt μ is defined as

$$\int_{\Omega} f(x) \mu(dx) = \inf_{s \geq f} \int_{\Omega} s(x) \mu(dx)$$

where the infimum is over simple functions. If the integral is finite, f is integrable.

- Extending to a real measurable function $f: \Omega \rightarrow \mathbb{R}$, we define $f(x) = f_+(x) - f_-(x)$, where

$$f_+(x) = \max\{0, f(x)\}$$

$$f_-(x) = -\min\{0, f(x)\}$$

so the integral of f is defined as

$$\int_{\Omega} f(x) \mu(dx) = \int_{\Omega} f_+(x) \mu(dx) - \int_{\Omega} f_-(x) \mu(dx)$$

and f is integrable if both f_+ and f_- are integrable.

- The Lebesgue integral satisfies the following properties:

↳ $f(x) = g(x)$ a.e. $\Rightarrow \int f d\mu = \int g d\mu$.

↳ Applying a convergence theorem on a seq. of integrable functions $f_n(x)$, we find the

pointwise limit $f(x) = \lim_{n \rightarrow \infty} f_n(x)$ is also an integrable function w/

$$\boxed{\int_{\Omega} f(x) \mu(dx) = \int_{\Omega} \lim_{n \rightarrow \infty} f_n(x) \mu(dx) = \lim_{n \rightarrow \infty} \int_{\Omega} f_n(x) \mu(dx)}$$

(Monotone convergence theorem): $f_n(x) \leq f_{n+1}(x)$ a.e. and $\sup_n \int_{\Omega} f_n(x) \mu(dx) < \infty$.

Dominated convergence theorem: \exists integrable function $g(x)$ s.t. $|f_n(x)| \leq g(x)$ a.e.

↳ The Lebesgue integral w/ counting measure μ_c is simply a sum

$$\boxed{\int_{\Omega} f(x) \mu_c(dx) = \sum_{n \in \mathbb{N}} f(n)}$$

↳ When (Ω, \mathcal{F}, P) is a prob. space, and $X: \Omega \rightarrow \mathbb{R}$ is a RV, we write the Lebesgue integral

$$\boxed{E[X] = \int_{\Omega} X(\omega) P(d\omega)}$$

For Personal Use Only -bkwk2

Radon-Nikodym derivatives

- Let μ and ν be two measures on a measurable space (Ω, \mathcal{F}) . ν is absolutely continuous wrt μ (denoted as $\nu \ll \mu$) if for any measurable set $A \in \mathcal{F}$

$$|\mu(A) = 0 \Rightarrow \nu(A) = 0|$$

- If $\mu \ll \nu$ and $\nu \ll \mu$, then μ and ν are said to be equivalent.

- Let μ and ν be two finite measures on a common measurable space (Ω, \mathcal{F}) . If $\nu \ll \mu$, then there exists a measurable function $f(x)$ s.t. for all measurable sets $A \in \mathcal{F}$,

$$|\nu(A) = \int_A f(x) \mu(dx)|$$

f is unique a.e. and is the Radon-Nikodym derivative of ν wrt μ , denoted as

$$f = \frac{d\nu}{d\mu} \Leftrightarrow |\nu(dx) = f(x) \mu(dx)|$$

- Let (Ω, \mathcal{F}, P) be a prob space, and $X: \Omega \rightarrow \mathbb{R}$ a RV w/ the law $\mu_X(A) = P(X^{-1}(A))$. If $\mu_1 \ll \mu_X / \mu_2 \ll \mu_X$, then X is a continuous RV, and its Radon-Nikodym derivative $\frac{d\mu_X}{d\mu_1} / \frac{d\mu_X}{d\mu_2}$ is the r.v.'s probability density function / probability mass function.

Product measure spaces.

- The product space $(\Omega, \mathcal{F}, \mu)$ of $(\Omega_1, \mathcal{F}_1, \mu_1)$ and $(\Omega_2, \mathcal{F}_2, \mu_2)$ is defined as follows:

↳ THE sample space $\Omega = \Omega_1 \times \Omega_2$ is all pairs of Ω_1 and Ω_2 .

↳ THE σ -algebra $\mathcal{F} = \mathcal{F}_1 \otimes \mathcal{F}_2$ is the smallest σ -algebra s.t. $\forall A \in \mathcal{F}_1, B \in \mathcal{F}_2, A \times B \in \mathcal{F}$.

↳ THE measure $\mu = \mu_1 \otimes \mu_2$ is a measure s.t. $\forall A \in \mathcal{F}_1, B \in \mathcal{F}_2, \mu(A \times B) = \mu_1(A) \mu_2(B)$

- THE Fubini-Tonelli thm. states for measure spaces (X, Σ_X, μ_X) and (Y, Σ_Y, μ_Y) , where μ_X and μ_Y are σ -finite measures, and f is a measurable function on the product space, we have

$$\int_X \int_Y f(x,y) |\mu_Y(dy) \mu_X(dx)} = \int_Y \int_X f(x,y) |\mu_X(dx) \mu_Y(dy)} = \int_{X \times Y} f(x,y) |\mu_X \otimes \mu_Y(d(x,y))$$

IF THE INTEGRAL IS FINITE, f IS INTEGRABLE AND WE HAVE

$$\int_X \int_Y f(x,y) \mu_Y(dy) \mu_X(dx) = \int_Y \int_X f(x,y) \mu_X(dx) \mu_Y(dy) = \int_{X \times Y} f(x,y) \mu_X \otimes \mu_Y(d(x,y))$$

Markov chains

Discrete Markov chains

- Let $\{X_n\}_{n \in \mathbb{N}}$ be a random process taking values in a finite set $S = \{1, \dots, k\}$. $\{X_n\}_{n \in \mathbb{N}}$ is a Markov chain if for any values $x_0, \dots, x_{n-1} \in S$ we have

$$P(X_n = x_n | X_0 = x_0, \dots, X_{n-1} = x_{n-1}) = P(X_n = x_n | X_{n-1} = x_{n-1})$$

i.e. the future is independent of the past, given the future

- If, in addition, the distribution $P(X_n = x_n | X_{n-1} = x_{n-1})$ is independent of n , it is a time-homogeneous Markov Chain. We define the transition matrix P as

$$P_{ij} = P(X_n = j | X_{n-1} = i) \quad \forall i, j \in S$$

i.e. Markov chains generate random sequences by looking only at the previously generated value

- The dist. of the first value is the initial dist. and is often written as a vector μ , where its i -th component μ_i is $\mu_i = P(X_0 = i)$. It then follows that

$$P(X_n = i) = (\mu P^n)_i$$

- If $\{X_n\}_{n \in \mathbb{N}}$ is a time-homogeneous Markov chain w/ transition matrix P , then π is a stationary dist. of $\{X_n\}_{n \in \mathbb{N}}$ if

$$\pi P = \pi$$

- A Markov chain is irreducible if for any $i, j \in S$, there is an $n \in \mathbb{N}$ s.t.

$$P(X_n = j | X_0 = i) > 0$$

i.e. any state is reachable from any state.

- A Markov chain is aperiodic if for any $i \in S$, there is an $n \in \mathbb{N}$ s.t.

$$P(X_n = i | X_0 = i) > 0 \quad \text{and} \quad P(X_{n+1} = i | X_0 = i) > 0$$

- If a Markov chain is irreducible and aperiodic, then it is regular ergodic. A regular ergodic Markov chain has a limiting dist. (unique stationary dist.) π which puts the mass on every element of the state space S .

- For a regular ergodic Markov chain $\{X_n\}_{n \in \mathbb{N}}$ w/ limiting dist. π , let $V_i(n) = \sum_{j=0}^{n-1} \mathbb{I}(X_j = i)$ be the no. of visits to state i up until time n $\forall i \in S$. The ergodic thm. states that

$$\frac{V_i(n)}{n} \rightarrow \pi_i \quad a.s.$$

i.e. the prop. of time spent in state i converges to its stationary prob. π_i as $n \rightarrow \infty$.

- We can thus evaluate the expected value $E_{\pi}(f(x))$ for a function $f: S \rightarrow \mathbb{R}$ on the state space S

$$E\left[\frac{1}{n} \sum_{j=0}^{n-1} f(X_j)\right] = \sum_{i \in S} f(i) E\left[\frac{V_i(n)}{n}\right] \rightarrow \sum_{i \in S} f(i) \pi_i = E_{\pi}(f(x))$$

For Personal Use Only -bkwk2

General state space Markov chains

- A Markov transition kernel on a measurable space (X, Σ_X) is a function $P: X \times \Sigma_X \rightarrow [0, 1]$ s.t.
 - ↳ For each fixed $x \in X$, $P(x, \cdot)$ is a prob. measure on (X, Σ_X)
 - ↳ For each fixed $A \in \Sigma_X$, $P(\cdot, A)$ is a non-negative measurable function on X .
- For any measurable sets $A_i \in \Sigma_X$, fix a starting pt. $x \in X$, we have :

$$P(X_1 \in A_1 | X_0 = x) = P(x, A_1)$$

$$P(X_1 \in A_1, X_2 \in A_2 | X_0 = x) = \int_{x \in A_1} P(x, dx_1) P(x_1, A_2)$$

$$P(X_1 \in A_1, X_2 \in A_2, \dots, X_n \in A_n | X_0 = x) = \int_{x \in A_1} P(x, dx_1) \int_{x_1 \in A_2} P(x_1, dx_2) \cdots \int_{x_{n-1} \in A_n} P(x_{n-1}, dx_n) P(x_n, A_n)$$

Given an initial dist. μ for X_0 , then we have

$$P(X_0 \in A_0, X_1 \in A_1, \dots, X_n \in A_n) = \int_{x_0 \in A_0} \mu(dx_0) \int_{x_1 \in A_1} P(x_0, dx_1) \cdots \int_{x_{n-1} \in A_{n-1}} P(x_{n-2}, dx_{n-1}) P(x_{n-1}, A_n)$$

* The integrals are well-defined by the measurability of the function $P(\cdot, A): X \rightarrow [0, 1]$

- A random process $\{X_n\}_{n \in \mathbb{N}}$ is a time-homogeneous Markov chain w/ transition kernel P if for any values $x_0, \dots, x_n \in X$ and $A \in \Sigma_X$ we have

$$P(X_n \in A | X_0 = x_0, \dots, X_{n-1} = x_{n-1}) = P(X_n \in A | X_{n-1} = x_{n-1}) = P(x_{n-1}, A)$$

- There are a few operations involving transition kernels. Let $f: X \rightarrow \mathbb{R}$ be a measurable function, $M: \Sigma_X \rightarrow \mathbb{R}^F$ a measure and P, Q transition kernels.

① PQ is a transition kernel, corr. to first transitioning w/ P , and then Q .

$$PQ(x, A) = \int P(x, dy) Q(y, A)$$

* We denote n transitions of P as $P^n = P \circ \dots \circ P$.

② $Pf: X \rightarrow \mathbb{R}$ is a function corr. to the expected value of f after a transition w/ P starting at $x \in X$.

$$Pf(x) = \int P(x, dy) f(y)$$

③ mP is a measure after one transition, defined as

$$mP(A) = \int P(x, A) M(dx)$$

* M is stationary for P if $mP = M$

- When the transition kernel P has a Radon-Nikodym derivative (density) wrt to the Lebesgue measure, we denote it as $p(x, y)$, where we have

$$P(x, dy) = p(x, y) dy \quad \Leftrightarrow \quad P(x, A) = \int_A p(x, y) dy$$

* Discrete Markov chains are just a special case of general state space Markov chains.

The transition matrix P corr. to the Radon-Nikodym derivative of P wrt M_0 .

For Personal Use Only -bkwk2

Marton chain Monte Carlo

Detailed balance

- A Marton chain w/ transition kernel P is reversible / satisfies detailed balance w.r.t. the measure π^* if for all $A, B \in \Sigma_X$,

$$\int_A P(x, y) \pi^*(dx) = \int_B P(y, A) \pi^*(dy) \quad \Leftrightarrow \quad p(x, y) \pi(x) = p(y, x) \pi(y)$$

where $\pi(x)$ is the R-N derivative (density) of measure π^* , so $\pi^*(dx) = \pi(x)dx$.

$$(\int_A (\int_B p(x, y) dy) \pi(x) dx = \int_B (\int_A p(y, x) dx) \pi(y) dy \rightarrow p(x, y) \pi(x) = p(y, x) \pi(y))$$

i.e. the prob. of observing transition $A \rightarrow B$ is equal to that of observing transition $B \rightarrow A$.

- If the transition kernel P satisfies no detailed balance eqns for the density π , then

$$\pi P = \pi$$

i.e. π is the invariant density of the Marton chain w/ transition kernel P .

Metropolis-Hastings (MH) algorithm

- Suppose there is a target prob. measure π^* from which we want to draw samples, we consider simulating a random process from a Marton chain w/ an invariant density π .

- Consider a transition operator $P(x, dy)$ of the form

$$P(x, dy) = p(x, y) dy + r(x) \delta_x(dy)$$

where δ_x is the dirac measure (1 if the argument is in x , 0 w/o)

- If we have $p(x, x) = 0$ and $\int_{\mathbb{R}^d} P(x, dy) = 1$, then we have

$$1 = \int_{\mathbb{R}^d} P(x, dy) = \int_{\mathbb{R}^d} p(x, y) dy + \int_{\mathbb{R}^d} r(x) \delta_x(dy) \rightarrow r(x) = 1 - \int_{\mathbb{R}^d} p(x, y) dy$$

which is the prob. that the chain remains at x after a transition.

- If $p(x, y)$ satisfies reversibility, then π is the invariant density of the transition kernel P

$$\int_{\mathbb{R}^d} P(x, A) \pi^*(dx) = \int_{\mathbb{R}^d} (\int_A p(x, dy)) \pi(x) dx$$

$$= \int_{\mathbb{R}^d} (\int_A p(x, y) dy + \int_A r(x) \delta_x(dy)) \pi(x) dx$$

$$= \int_A \int_{\mathbb{R}^d} p(x, y) \pi(x) dy dx + \int_{\mathbb{R}^d} r(x) \delta_x(A) \pi(x) dx$$

$$= \int_A \int_{\mathbb{R}^d} p(y, x) \pi(x) dy dx + \int_A r(x) \pi(x) dx = \int_A \pi(y) dy = \pi^*(A)$$

detailed balance

- Consider a proposal density $q(x, y)$ s.t. $\int q(x, y) dy = 1$.

If $x \rightarrow y$ transitions are more freq. than $y \rightarrow x$ transitions, we have $q(x, y) \pi(x) > q(y, x) \pi(y)$

We rebalance both sides by introducing $\alpha(x, y) < 1$, $\alpha(x, y) q(x, y) \pi(x) = q(y, x) \pi(y)$.

$$\alpha(x, y) = \frac{q(y, x) \pi(y)}{q(x, y) \pi(x)}$$

On the contrary, if we have $q(x, y) \pi(x) < q(y, x) \pi(y)$, set $\alpha(x, y) = 1$, so for $x \neq y$,

$$p(x, y) = \alpha(x, y) q(x, y), \quad \text{where } \alpha(x, y) = \min \left\{ \frac{\pi(y) q(y, x)}{\pi(x) q(x, y)}, 1 \right\}$$

For Personal Use Only -bkwk2

- The overall transition operator P is thus given by

$$P(x, dy) = \alpha(x, y) q(x, y) dy + [1 - \int_{\mathbb{R}^d} \alpha(x, y) q(x, y) dy] \delta_x(dy) \quad \text{or}$$

$$P(x, A) = \int_A \alpha(x, y) q(x, y) dy + [1 - \int_{\mathbb{R}^d} \alpha(x, y) q(x, y) dy] \delta_x(A)$$

which is reversible by construction, therefore has π as its invariant density.

* Note that the normalizing const. of the target density $\pi(x) = \int \phi(x)$ is not req.

$$\alpha(x, y) = \min \left\{ \frac{\phi(y)}{\phi(x)} \frac{q(y, x)}{q(x, y)}, 1 \right\}$$

If the proposal density is symmetric, i.e. $q(x, y) = q(y, x)$,

$$\alpha(x, y) = \min \left\{ \frac{\phi(y)}{\phi(x)}, 1 \right\}.$$

- Most commonly, $q(x, y)$ is chosen as a Normal dist., i.e. $q(x, y) \propto \exp(-\frac{\|x-y\|^2}{2\sigma^2})$.

However, it is not trivial to optimally choose σ^2 .

↳ σ^2 too small: Markov chain move very slowly across state space

↳ σ^2 too big: A proposal is rarely accepted.

* Both cases yield high correlation samples \rightarrow bad mixing.

- The MH algorithm has a few drawbacks:

↳ It is ineffective for sampling multimodal dist. — especially if the target dist. π has two highly conc. and distant modes (chain stuck in one of them).

↳ Reversible Markov chains are likely to make steps back, and therefore explore the state space in a slow, diffusive fashion.

Pseudo-code for the MH algorithm.

- Given the target density π , a proposal density q , the first state x_0 , and the no. of iterations N ,

for $i = 1, \dots, N$ do

$y \sim Q(x_{i-1}, \cdot)$

$u \sim U[0, 1]$

if $u < \alpha(x_{i-1}, y)$ then

$x_i = y$

N iterations

make a proposal using $Q(x_{i-1}, \cdot)$

Accept w/ prob. $\alpha(x_{i-1}, y)$, set $x_i = y$

else

$x_i = x_{i-1}$

Rej. w/ prob. $1 - \alpha(x_{i-1}, y)$, set $x_i = x_{i-1}$

end if

end for

return x_0, \dots, x_N

For Personal Use Only -bkwk2

Mixture transition kernels

- The mixture transition kernel is composed of the mixture

$$\gamma P_1(x, dy) + (1-\gamma) P_2(x, dy), \quad \gamma \in [0, 1]$$

where transition kernels $P_1(x, dy)$ and $P_2(x, dy)$ both have invariant density $\pi(x)$.

- The mixture transition kernel has invariant density $\pi(x)$

$$\int_{\mathbb{R}^d} [\gamma P_1(x, dy) + (1-\gamma) P_2(x, dy)] \pi^*(dx) = \pi^*(dy)$$

$$(LHS = \gamma \int_{\mathbb{R}^d} P_1(x, dy) \pi(x) dx + (1-\gamma) \int_{\mathbb{R}^d} P_2(x, dy) \pi(x) dx = \gamma \pi^*(dy) + (1-\gamma) \pi^*(dy) = \pi^*(dy))$$

- The mixture transition kernel is useful for drawing samples from multimodal target densities.

Product transition kernels

- Consider breaking up a vector up into sub-blocks, $x = (x_1, x_2)$ where $x \in \mathbb{R}^{d_1+d_2}$, $x_1 \in \mathbb{R}^{d_1}$, $x_2 \in \mathbb{R}^{d_2}$
- The product transition kernel is given by

$$P_1(x_1, dy_1 | x_2) P_2(x_2, dy_2 | y_1)$$

where conditional transition kernels $P_1(x_1, dy_1 | x_2)$ and $P_2(x_2, dy_2 | y_1)$ have invariant densities

$\pi_{1|2}(\cdot | x_2)$ for fixed x_2 , $\pi_{2|1}(\cdot | y_1)$ for fixed y_1 , respectively

$$(i.e. \pi_{1|2}(dy_1 | x_2) = \int P_1(x_1, dy_1 | x_2) \pi_{1|2}(x_1 | x_2) dx_1; \pi_{2|1}(dy_2 | y_1) = \int P_2(x_2, dy_2 | y_1) \pi_{2|1}(x_2 | y_1) dx_2)$$

- The product transition kernel has invariant density $\pi(x_1, x_2) = \pi_{1|2}(x_1 | x_2) \pi_{2|1}(x_2) = \pi_{2|1}(x_2 | y_1) \pi_1(y_1)$

$$\int_{\mathbb{R}^{d_2}} \int_{\mathbb{R}^{d_1}} P_1(x_1, dy_1 | x_2) P_2(x_2, dy_2 | y_1) \pi^*(dx_1, dx_2) = \pi^*(dy_1, dy_2)$$

$$(LHS = \int_{\mathbb{R}^{d_2}} \int_{\mathbb{R}^{d_1}} P_1(x_1, dy_1 | x_2) P_2(x_2, dy_2 | y_1) \pi(x_1, x_2) dx_1 dx_2$$

$$= \int_{\mathbb{R}^{d_2}} P_2(x_2, dy_2 | y_1) \left[\int_{\mathbb{R}^{d_1}} P_1(x_1, dy_1 | x_2) \pi_{1|2}(x_1 | x_2) dx_1 \right] \pi_2(x_2) dx_2$$

$$= \int_{\mathbb{R}^{d_2}} P_2(x_2, dy_2 | y_1) \pi_{1|2}(dy_1 | x_2) \pi_2(x_2) dx_2 = \int_{\mathbb{R}^{d_2}} P_2(x_2, dy_2 | y_1) \frac{\pi_{2|1}(x_2 | y_1) \pi_1(y_1)}{\pi_2(x_2)} \pi_2(x_2) dx_2$$

$$= \pi_1(y_1) \int_{\mathbb{R}^{d_2}} P_2(x_2, dy_2 | y_1) \pi_{2|1}(x_2 | y_1) dx_2 = \pi_1(y_1) \pi_{2|1}^*(y_2 | y_1) = \pi^*(dy_1, dy_2)$$

Gibbs sampling (GS) algorithm

- In GS, for each component j of x , x_j , in turn, we sample a new value from the conditional distribution of x_j given all other variables x_{-j} .

$$x_j \sim p(x_j | x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_d)$$

This is equivalent to using the product transition kernel $\prod_{j=1}^D P_j(x_j, dy_j | x_{-j})$ and using MH w/ exact conditionals as proposals.

- For $D=2$, the transition $(x_1, x_2) \rightarrow (y_1, y_2)$ consists of $(x_1, x_2) \xrightarrow{[1]} (y_1, x_2) \xrightarrow{[2]} (y_1, y_2)$

$$\hookrightarrow \text{step [1]: } q([x_1, x_2], [y_1, x_2]) = \pi(y_1 | x_2); q([y_1, x_2], [x_1, x_2]) = \pi(x_1 | x_2)$$

$$\therefore \frac{\pi(y_1, x_2) q([y_1, x_2], [x_1, x_2])}{\pi(x_1, x_2) q([x_1, x_2], [y_1, x_2])} = \frac{\pi(y_1 | x_2) \pi(x_2)}{\pi(x_1 | x_2) \pi(x_2) \pi(y_1 | x_2)} = 1 \rightarrow q([x_1, x_2], [y_1, x_2]) = 1$$

$$\hookrightarrow \text{step [2]: } q([y_1, x_2], [y_1, y_2]) = \pi(y_2 | y_1); q([y_1, y_2], [y_1, x_2]) = \pi(x_2 | y_1)$$

$$\therefore \frac{\pi(y_1, y_2) q([y_1, y_2], [y_1, x_2])}{\pi(y_1, x_2) q([y_1, x_2], [y_1, y_2])} = \frac{\pi(y_1 | y_2) \pi(y_2)}{\pi(y_1 | x_2) \pi(x_2) \pi(y_2 | y_1)} = 1 \rightarrow q([y_1, x_2], [y_1, y_2]) = 1$$

For Personal Use Only -bkwk2

Pseudo-code for the GS algorithm.

- Given the target density π , the first state x_0 and no. of iterations N ,

for $i=1, \dots, N$ do

 for $j=1, \dots, D$ do

$$x_i^j \sim \pi(\cdot | x_i^{j-1}, x_i^{j+1}, \dots, x_N)$$

 end for

end for

return x_0, \dots, x_N

Data augmentation

- Consider a target density $\pi(\theta)$ where $\theta \in \mathbb{R}^p$. Now augment the model w/ $\phi \in \mathbb{R}^q$,

so the augmented target density is $\pi(\theta, \phi)$.

- The desired density can be recovered as

$$\pi(\theta) = \int \pi(\theta, \phi) d\phi$$

→ We consider a Markov chain w/ invariant density $\pi(\theta, \phi)$ and draw samples $\theta^{(n)}, \phi^{(n)} \sim \pi(\theta, \phi)$

For each $\phi^{(n)}$, it follows that $\pi(\theta^{(n)}) = \int \pi(\theta^{(n)}, \phi) d\phi$ - each $\theta^{(n)}$ marginally dist. w/ $\pi(\theta)$.

- We select ϕ s.t. the exact conditionals $\pi(\phi|\theta)$ and $\pi(\theta|\phi)$ can be sampled directly.

Hilbert spaces and Sobolev spaces

Vector, metric, Banach and Hilbert spaces

Vector space

- A vector space V is a collection of vectors $x, y, z \in V$ that satisfy the following axioms.
- ↳ Addition commutative $x+y = y+x \in V$
- ↳ Addition associative $x+(y+z) = (x+y)+z \in V$
- ↳ Zero vector exists $\underline{x+0=x}$
- ↳ Additive inverse exists $x+(-x) = \underline{0}$
- ↳ Scalar multiplication closed $\alpha x \in V, \alpha \in F$
- ↳ Scalar multiplication distributive $\alpha(x+y) = \alpha x + \alpha y \in V, \alpha \in F$
- ↳ Scalar multiplication associative $\alpha(\beta x) = \beta(\alpha x), \alpha, \beta \in F$.
- ↳ Zero and unit scalar $0x = \underline{0} \in V, 1x = \underline{x}$.
- Vector space is a group under addition.

Inner product space

- An inner product space is a vector space V over the field F together with an inner product, which is the mapping $\langle \cdot, \cdot \rangle : V \times V \rightarrow F$ that satisfies $\forall x, y, z \in V, \forall \alpha, \beta \in F$.
- ↳ Conjugate symmetry $\langle x, y \rangle = \overline{\langle y, x \rangle} \rightarrow \langle x, y \rangle \in R$
- ↳ Linearity $\langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle$
- ↳ Positive definiteness $\langle x, x \rangle > 0 \text{ if } x \neq 0 \quad (\langle x, x \rangle = 0 \text{ iff } x = 0)$
- For vectors $a, b \in C^n$, the inner product is defined as

$$\langle a, b \rangle = \sum_{i=1}^n a_i^* b_i$$

For functions f, g of variable on interval $x \in [a, b] \subset R$, the inner product is defined as

$$\langle f, g \rangle = \int_a^b f(x)g(x) dx$$

Metric space

- A metric space is a vector space V that has a metric function, which is a mapping $\rho : V \times V \rightarrow R$ that satisfies $\forall x, y, z \in V$
- ↳ zero self-distance $\rho(x, x) = 0 \text{ iff } x = y$
- ↳ positivity $\rho(x, y) > 0 \text{ iff } x \neq y$
- ↳ symmetry $\rho(x, y) = \rho(y, x)$
- ↳ triangle inequality $\rho(x, z) \leq \rho(x, y) + \rho(y, z)$

For Personal Use Only - bkwk2

Normed spaces

- A normed space is a vector space V over the field \mathbb{F} together w/ a norm, which is a mapping $\|\cdot\|: V \rightarrow \mathbb{R}$ that satisfies $\forall u, v \in V, \forall \lambda \in \mathbb{F}$
- ↳ Non-negativity $\|u\| \geq 0$
- ↳ Positive-definiteness $\|u\| = 0 \iff u = 0$
- ↳ Absolute homogeneity $\|\lambda u\| = |\lambda| \|u\|$
- ↳ Triangle inequality $\|u+v\| \leq \|u\| + \|v\|$

- For inner product spaces V , the norm can be induced by the inner product $\langle \cdot, \cdot \rangle$,

$$\|u\| = \sqrt{\langle u, u \rangle}$$

For any two elements $x, y \in V$, we have the following (in)equalities:

↳ Parallelogram law: $\|x+y\|^2 + \|x-y\|^2 = 2(\|x\|^2 + \|y\|^2)$

$$(\langle x+y, x+y \rangle + \langle x-y, x-y \rangle = \langle x, x \rangle + 2\langle x, y \rangle + \langle y, y \rangle + \langle x, x \rangle - 2\langle x, y \rangle + \langle y, y \rangle)$$

↳ Cauchy-Schwarz inequality $|\langle x, y \rangle| \leq \|x\| \|y\|$

(Let $z = x - \lambda y$ for some scalar λ .

$$\|z\|^2 = \|x - \lambda y\|^2 = \|x\|^2 - 2\lambda \langle x, y \rangle + \lambda^2 \|y\|^2 \geq 0$$

→ we have $A\lambda^2 + B\lambda + C \geq 0$. where $A = \|y\|^2$, $B = -2\langle x, y \rangle$, $C = \|x\|^2$

since $A\lambda^2 + B\lambda + C$ is non-negative, we have $\Delta = B^2 - 4AC \leq 0$,

$$(-2\langle x, y \rangle)^2 - 4\|x\|^2\|y\|^2 \geq 0 \rightarrow |\langle x, y \rangle| \geq (\|x\| \|y\|)$$

* Lp, L-p norms are defined w/o reference to an inner product → norm is more general.

* A normed space is also a metric space. w/ $d(x, y) = \|x - y\|$

Orthogonality and orthonormality

- For a set of vectors $\{v_1, v_2, \dots, v_n\}$ in an inner product space to be orthogonal, each pair of distinct vectors are orthogonal,

$$\langle v_i, v_j \rangle = 0 \quad \forall i \neq j$$

If each vector is a unit vector, the set of vectors is orthonormal.

$$\|v_i\| = 1 \quad \forall i$$

- A set of vectors v_1, v_2, \dots, v_n is linearly independent iff

$$\sum_{i=1}^n c_i v_i = 0 \Rightarrow c_i = 0 \quad \forall i$$

Any set of orthonormal vectors e_1, e_2, \dots, e_n is linearly independent, thus can serve as a basis for a finite dim. inner product space V . Every vector $v \in V$ can be written

$$v = \sum_{i=1}^n \alpha_i e_i$$

where $\alpha_1, \alpha_2, \dots, \alpha_n$ are the coordinates of v .

For Personal Use Only -bkwk2

Completeness

- A sequence of vectors $\{x_1, x_2, x_3, \dots\}$ is a Cauchy sequence of vectors if for any positive $\epsilon > 0$, there exists a no. N s.t.

$$\|x_m - x_n\| < \epsilon \quad \forall m, n > N$$

i.e. the seq. of terms x_m, x_n get closer and closer in an unlimited manner as $m, n \rightarrow \infty$.

- An infinite sequence of vectors $\{x_1, x_2, \dots\}$ is said to be convergent if there exists an element $x \in V$ s.t.

$$\|x_n - x\| \rightarrow 0$$

- If every Cauchy sequence in a vector space V converges to a limit that is also in V , then the space V is complete.

* All finite dim. normed spaces are complete (not necessarily for infinite dim.).

Banach spaces

- A Banach space is a normed space that is complete

- All finite dim. normed spaces are complete \rightarrow Banach spaces.

Not all infinite dim. normed spaces are complete (some are, e.g. L^p, L^∞ spaces).

- The l^p space consists of all elements $x = (x_1, x_2, \dots)$ s.t. its l^p norm is finite,

$$\|x\|_{l^p} = \left(\sum_{i=1}^{\infty} |x_i|^p \right)^{1/p} < \infty \quad \text{for } 1 \leq p < \infty$$

For $p=\infty$,

$$\|x\|_{l^\infty} = \sup_i |x_i| < \infty$$

- The L^p space consists of measurable functions f defined on the measure space (X, Σ, μ) s.t. its L^p norm is finite,

$$\|f\|_{L^p} = \left(\int_X |f(x)|^p d\mu \right)^{1/p} < \infty \quad \text{for } 1 \leq p < \infty$$

For $p=\infty$,

$$\|f\|_{L^\infty} = \text{ess sup}_{x \in X} |f(x)| < \infty$$

- The l^p, L^p spaces are normed and complete \rightarrow they are Banach spaces

* Only for $p=2$, l^p, L^p spaces have an inner product (norm induced by inner product)

Pre-Hilbert space

- A pre-Hilbert space is an inner product space that is incomplete.

- e.g.: Consider the seq. of functions $\{f_n\}$ defined on the unit line.

$$f_n(x) = \begin{cases} 1 & x \in [0, 1/2] \\ 1 - 2n(x - 1/2) & x \in [1/2, 1/(1+n)] \\ 0 & x \in [1/(1+n), 1] \end{cases}$$

$$\text{Note that } \|f_n - f_m\|_{L_2} = \sqrt{\int_0^1 (f_n(x) - f_m(x))^2 dx} = \left(1 + \frac{m}{n}\right) \sqrt{\int_0^1 f_n^2(x) dx}, \lim_{m, n \rightarrow \infty} \|f_n - f_m\|_{L_2} = 0$$

$\Rightarrow \{f_n\}$ is a Cauchy seq. that converges to $f(x) = \begin{cases} 1 & x \in [0, 1/2] \\ 0 & x \in [1/2, 1] \end{cases}$

$\{f_n\}$ is continuous, f is not continuous \rightarrow space is incomplete \rightarrow pre-Hilbert space.

For Personal Use Only -bkwk2

Hilbert space

- A Hilbert space is a Banach space endowed w/ an inner product.
- Let $\{e_i\}$ be an infinite set of orthonormal vectors in Hilbert space H .

For any function f (a pt in infinite dim. space H) denote the coefficients $c_i = \langle e_i, f \rangle$.

↳ Bessel's inequality:

$$\sum_{i=1}^{\infty} c_i^2 \leq \|f\|^2$$

(Consider finite sum from first n terms in the infinite seq., $f_n = \sum_{i=1}^n c_i e_i$.

$$\|f_n\|^2 = \langle f_n, f_n \rangle = \sum_{i=1}^n \sum_{j=1}^n c_i c_j \langle e_i, e_j \rangle = \sum_{i=1}^n \sum_{j=1}^n c_i c_j \delta_{i,j} = \sum_{i=1}^n c_i^2$$

$$\langle f, f_n \rangle = \sum_{i=1}^{\infty} \sum_{j=1}^n c_i c_j \langle e_i, e_j \rangle = \sum_{i=1}^{\infty} \sum_{j=1}^n c_i c_j \delta_{i,j} = \sum_{i=1}^n c_i^2$$

Cauchy-Schwarz: $|\langle f, f_n \rangle|^2 \leq \|f\|^2 \|f_n\|^2 \rightarrow \sum_{i=1}^n c_i^2 \leq \|f\|^2$. Holds for $n=\infty$.)

* In l^2 / L^2 spaces, $\|f\| < \infty$, by Bessel's inequality. $\sum_{i=1}^{\infty} c_i^2 < \infty$.

↳ Parseval's identity:

$$\{e_i\} \text{ is complete} \Leftrightarrow \sum_{i=1}^{\infty} c_i^2 = \|f\|^2$$

* We can use Parseval's identity to check that an infinite basis $\{e_i\}$ is complete.

- l^2, L^p spaces have an inner product for $p=2 \rightarrow l^2, L^2$ spaces are Hilbert spaces.

- A sequence of functions in L^2 space converge if for any small $\epsilon > 0$, there

exists a no. N s.t. $\|f_n - f\|_{L_2} < \epsilon \ \forall n > N$. As $\|f\|_{L_2} = \sqrt{\int_a^b |f(x)|^2 dx}$, we have

$$\|f_n - f\|_{L_2} = \sqrt{\int_a^b |f_n(x) - f(x)|^2 dx} \rightarrow 0$$

i.e. convergence in the L_2 norm means convergence a.e. $\rightarrow f, f_n$ are equivalence classes.

* We only req. equivalence a.e. as sets w/ zero measure don't contribute to the integral.

$\rightarrow L^2$ space is a space of equivalence functions, i.e. each equivalence class is grouped as the same pt. in the L^2 space.

- Consider a set of orthonormal functions $\{\phi_k\}$ that are square integrable ($\int k(x)^2 dx < \infty$).

Any function $f \in L^2$ can be written as a Fourier series wrt $\{\phi_k\}$.

$$f = \sum_{k=1}^{\infty} c_k \phi_k$$

where $c_k = \langle f, \phi_k \rangle$ are the Fourier coefficients of f .

By Bessel's inequality, $\sum_{k=1}^{\infty} c_k^2 \leq \|f\|_{L_2}^2$, and since $f \in L^2$, $\|f\|_{L_2} < \infty$, therefore.

$$\sum_{k=1}^{\infty} c_k^2 < \infty$$

\therefore The seq. of Fourier coeff $\{c_k\}$ is an element of the l^2 space.

\rightarrow The elements $f \in L^2$ and $c = (c_1, c_2, \dots) \in l^2$ are connected via the Fourier coeff $c_k = \langle f, \phi_k \rangle$.

The l^2, L^2 spaces are isomorphic, and l^2 is a coordinate system for L^2 .

For Personal Use Only -bkwk2

reproducing kernel Hilbert space (RKHS)

- A RKHS is a Hilbert space of functions in which point evaluation is a continuous linear functional - if functions are close in norm, they are also close pointwise, i.e.

$$\|f - g\|_{L_2} \rightarrow 0 \Rightarrow |f(x) - g(x)| \rightarrow 0 \quad \forall x.$$

- A RKHS defines a reproducing kernel K_x , which is symmetric and positive definite, s.t.

$$f(x) = \langle f, K_x \rangle$$

Conversely, Moore-Aronszajn thm. states that a reproducing kernel K_x uniquely defines a RKHS.

- consider a dataset $(x_1, y_1), \dots, (x_N, y_N) \in X \times \mathbb{R}$ representing $y_i = f(x_i)$.

By calculus of variations, the regularised empirical error $\sum_{n=1}^N |f(x_n) - y_n|^2 + \lambda \|f\|^2$ is minimized by functions of the form

$$f(\cdot) = \sum_{n=1}^N \alpha_n K(\cdot, x_n)$$

where $\alpha = (\alpha_1, \dots, \alpha_N)$ are given by the sol'n of the linear system

$$(K + (N\lambda) I_N) \alpha = y$$

where λ is a scalar, K is a $N \times N$ matrix w/ elements $K(x_i, x_j)$,

and y is a $N \times 1$ vector w/ elements y_i

→ This enables function approximation via the form $f(\cdot) = \sum_{n=1}^N \alpha_n K(\cdot, x_n)$

↳ Proof: $\|f\|^2 = \langle f, f \rangle = \sum_i \sum_j \alpha_i \alpha_j \langle K(\cdot, x_i), K(\cdot, x_j) \rangle = \sum_i \sum_j \alpha_i \alpha_j K(x_i, x_j) = \alpha^T K \alpha$

$$\therefore \frac{\partial}{\partial \alpha} \frac{1}{N} (K \alpha - y)^T (K \alpha - y) + \lambda \alpha^T K \alpha = 0 \rightarrow \frac{1}{N} K (K \alpha - y) + \lambda K \alpha = 0$$

$$\Rightarrow \alpha^* = \frac{\alpha^T K^T}{\alpha^T K^T K} \frac{1}{N} (K \alpha - y)^T (K \alpha - y) + \lambda \|f\|^2 = (K + (N\lambda) I_N)^{-1} y$$

Sobolev spaces and function approximation

Weak derivative

- consider the function $f: [0,1] \rightarrow \mathbb{R}$. The function $h: [0,1] \rightarrow \mathbb{R}$ is a weak derivative of f if for every differentiable function $g: [0,1] \rightarrow \mathbb{R}$ w/ $g(0) = g(1) = 0$ it holds that

$$\int_0^1 f'(x) g'(x) dx = - \int_0^1 h(x) g(x) dx$$

- The notion of weak derivatives can be generalized to multiple dim over more general domains,

and higher derivatives follow straightforwardly (not important - just know they exist + are well-defined)

- e.g.: consider the function $f: [0,2] \rightarrow [0,1]$, $f(x) = \begin{cases} x & x \in [0,1] \\ 1 & x \in [1,2] \end{cases}$

for any differentiable test function $\phi: [0,2] \rightarrow \mathbb{R}$, w/ $\phi(0) = \phi(2) = 0$, then

$$-\int_0^2 x \phi'(x) dx = -[x \phi(x)]_0^2 + \int_0^2 \phi(x) dx = -\phi(1) + \int_0^2 \phi(x) dx ; -\int_0^2 \phi'(x) dx = -[\phi(x)]_0^2 = \phi(1)$$

$$\therefore -\int_0^2 f(x) \phi'(x) dx = \int_0^2 \phi(x) dx = \int_0^2 \mathbb{I}_{x \in [0,1]} dx \rightarrow Df(x) = \mathbb{I}_{x \in [0,1]}$$

For Personal Use Only -bkwk2

Sobolev spaces

- Denote a space of smooth continuous functions of degree k (has k derivatives), defined on a domain which is an open set on the real line SCR as $C^k(\Omega)$.
The space $C^k(\Omega)$ is incomplete \rightarrow enlarge to a Lebesgue space, L^p .
- Consider a function $f \in L^p(\Omega)$ defined on SCR. Let f have weak derivatives up to k th degree which also belong to $L^p(\Omega)$. A Sobolev space $W_p^k(\Omega)$ is defined by the set of functions

$$\rightarrow W_p^k(\Omega) = \{f \in L^p(\Omega) : D^\alpha f \in L^p(\Omega), 0 \leq |\alpha| \leq k\}$$

w/ the norm defined as follows:

$$\|f\|_{W_p^k} = \left(\int_{\Omega} |f(x)|^p dx + \int_{\Omega} |Df(x)|^p dx + \dots + \int_{\Omega} |D^k f(x)|^p dx \right)^{1/p}$$

- For the case $p=2$, $W_2^k(\Omega)$ is a Hilbert space w/ the inner product

$$\langle f, g \rangle_{W_2^k(\Omega)} = \sum_{|\alpha|=0}^k \langle D^\alpha f, D^\alpha g \rangle$$

* Both the norm and inner product encode aspects of the smoothness of the function classes.

Clarify for
new case,

we which
 $D^\alpha f(\cdot)$

Function approximation by polynomial expansion

- Consider functions from the set $C_2[-\pi, \pi] = C[-\pi, \pi] \cap L^2[-\pi, \pi]$ (continuous + square integrable)

We can represent functions from this set using Fourier expansions

$$f(x) = \sum_{k=-\infty}^{\infty} c_k e^{ikx}, \quad c_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-ikx} dx$$

Note the L_2 norm is given by $\|f\|_{L_2}^2 = \sum_{k=-\infty}^{\infty} c_k^2 < \infty$.

$$\|f(x)\|_{L_2}^2 = \int_{-\pi}^{\pi} \left(\sum_i c_i \phi_i(x) \right) \left(\sum_j c_j \phi_j(x) \right) dx = \sum_i \sum_j c_i c_j \int_{-\pi}^{\pi} \phi_i(x) \phi_j(x) dx = \sum_i \sum_j c_i c_j \delta_{ij} = \sum_{k=-\infty}^{\infty} c_k^2$$

- Now consider the Sobolev space W_2^k , defined here as $W_2^k = \{f \in C_2[-\pi, \pi] : \|D^\alpha f\|_{L_2}^2 < \infty\}$.

\rightarrow Here we define the norm as $\|f\|_{W_2^k}^2 = \|D^\alpha f\|_{L_2}^2$, which is given by $\|f\|_{L_2}^2 = \sum_{k=-\infty}^{\infty} k^{2s} c_k^2 < \infty$.
 $(\|f(x)\|_{L_2}^2 = \|D^\alpha f\|_{L_2}^2 = \left\| \sum_k (ik)^s c_k e^{ikx} \right\|_{L_2}^2 = \sum_k (ik)^{2s} c_k^2 = \sum_k k^{2s} c_k^2)$ ↑s → high order terms
in Fourier series have less contribution.

* For the norm to be bounded, c_k need to decay at a rate which increases w/ s .

- Consider an increasing space H_n as the set of trigonometric polynomials of degree n . The optimal approx. is

$$f_n(x) = \sum_{k=-n}^n c_k e^{ikx}$$

For $f \in W_2^k$, the approximation error is given by $E_n[f] < \frac{1}{n^{2s}} \|f\|_{W_2^k}^2$

$$(E_n[f])^2 = \|f - f_n\|_{L_2}^2 = \left\| \sum_{k=n+1}^{\infty} c_k e^{ikx} \right\|_{L_2}^2 = \sum_{k=n+1}^{\infty} c_k^{2s} \frac{k^{2s}}{k^{2s}} < \frac{1}{n^{2s}} \sum_{k=n+1}^{\infty} k^{2s} c_k^2 < \frac{1}{n^{2s}} \sum_{k=1}^{\infty} k^{2s} c_k^2 = \frac{1}{n^{2s}} \|f\|_{W_2^k}^2$$

* The rate of convergence is faster and impact of increasing model complexity n is greater as the smoothness of the function being approx. increases

- For a function in a $d=2$ dim. domain, we have $f(x) = \sum_{k,m=1}^n c_{k,m} e^{ikx_1 + imx_2}$

$$\|f\|_{W_2^k}^2 = \|D^\alpha f\|_{L_2}^2 = \|D_x^\alpha f_1\|_{L_2}^2 + \|D_x^\alpha f_2\|_{L_2}^2 = \sum_{k,m=1}^n k^{2s} c_{k,m}^2 + \sum_{k,m=1}^n m^{2s} c_{k,m}^2 = \sum_{k,m=1}^n (k^2 + m^2) c_{k,m}^2 \rightarrow E_n[f] < \frac{1}{n^{2s}} \|f\|_{W_2^k}^2$$

For a function in a d dim. domain, we have $f(x) = \sum_{i_1, \dots, i_d=1}^n c_{i_1, \dots, i_d} e^{i_1 x_1 + \dots + i_d x_d}$

$$\|f\|_{W_2^k}^2 = \sum_{i_1, \dots, i_d=1}^n (i_1^2 + \dots + i_d^2) c_{i_1, \dots, i_d}^2 \rightarrow E_n[f] < \frac{1}{d n^{2s/d}} \|f\|_{W_2^k}^2$$

For Personal Use Only -bkwk2

Gaussian measure in Hilbert space

Lebesgue measure in infinite dimensions

- The Lebesgue measure m_L is a generalization of length. In \mathbb{R} , we have $m_L([a,b]) = b-a$.
- The Lebesgue measure m_L has the following properties:
 - ↪ FINITENESS and monotonicity: If $A \subset B \subset \mathbb{R}$, then $0 \leq m_L(A) \leq m_L(B) \leq \infty$
 - ↪ Translation invariance: If $A \subset \mathbb{R}$, $x_0 \in \mathbb{R}$, $m_L(A+x_0) = m_L(A)$, where $A+x_0 = \{x+x_0 : x \in A\}$.
- In \mathbb{R}^D , we have $m_L([a_1, b_1] \times \dots \times [a_D, b_D]) = \prod_{i=1}^D (b_i - a_i)$.

Consider the case $D=\infty$. For a Hilbert space H , there is a countable orthonormal set acting as a basis in H , i.e. $\{e_i : i=1, 2, \dots\}$ w/ $\langle e_i, e_j \rangle = \delta_{ij}$.

We have $\|e_i - e_j\|^2 = \|e_i\|^2 - 2 \langle e_i, e_j \rangle + \|e_j\|^2 = 1 + 1 = 2$ if $i \neq j$.

Denoting a ball centered at c w/ radius r as $B(c, r)$, place $\{B(e_i, \frac{r}{2})\}_{i \in \mathbb{N}}$ in $B(0, 2)$.

→ we have $B(e_i, \frac{r}{2}) \cap B(e_j, \frac{r}{2}) = \emptyset$ and $\bigcup_{i \in \mathbb{N}} B(e_i, \frac{r}{2}) \subset B(0, 2)$. m_L is σ -finite

Monotonicity + countable additivity of m_L : $\sum_{i \in \mathbb{N}} m_L(B(e_i, \frac{r}{2})) \leq m_L(B(0, 2)) < \infty$. [1]

Translation invariance of m_L : $B(e_i, \frac{r}{2}) = B(e_j, \frac{r}{2}) \rightarrow \sum_{i \in \mathbb{N}} m_L(B(e_i, \frac{r}{2})) = \text{const} \times \sum_{i \in \mathbb{N}} = \infty$. [2]

⇒ [1], [2] contradict, so there is no Lebesgue measure m_L in infinite dim. Hilbert space.

Gaussian measure in infinite dimensions

- Define the infinite dimensional product measure on \mathbb{R}^∞ as

$$m = \prod_{n=1}^{\infty} m_n$$

where each m_n is the standard Gaussian in \mathbb{R} , $m_n(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x^2) dx$

- Consider the function f_λ

$$f_\lambda(x) = \exp\left(-\sum_{k=1}^{\infty} \frac{1}{2} a_k x_k^2\right)$$

which is only well-defined (finite) if $\sum_{k=1}^{\infty} a_k x_k^2 < \infty$.

We define the Hilbert space $\ell_{2,a}$ w/ inner product $\langle x, y \rangle = \sum_{k=1}^{\infty} a_k x_k y_k$.

$$x \in \ell_{2,a} = \{x \in \mathbb{R}^\infty : \sum_{k=1}^{\infty} a_k x_k^2 = \|x\|_{2,a}^2 < \infty\}$$

Taking the Lebesgue integral wrt infinite dim. standard Gaussian product measure,

$$\int_{\mathbb{R}^\infty} f_\lambda d\mu = \int_{\mathbb{R}^\infty} f_\lambda \prod_{k=1}^{\infty} d\mu_k = \prod_{k=1}^{\infty} \int_{\mathbb{R}} \exp\left(-\frac{1}{2}(1+\lambda a_k)x_k^2\right) dx_k = \prod_{k=1}^{\infty} \sigma_k \int_{-\infty}^{\infty} N(x_k | 0, \sigma_k^2) dx_k = \prod_{k=1}^{\infty} \sigma_k$$

where $\sigma_k = \sqrt{1+\lambda a_k}$, so we have $\int_{\mathbb{R}^\infty} f_\lambda d\mu = \prod_{k=1}^{\infty} \frac{1}{\sqrt{1+\lambda a_k}}$

* If $\sum_{k=1}^{\infty} a_k < \infty$, then $\prod_{k=1}^{\infty} (1+\lambda a_k) < \infty$.

As $\lambda \rightarrow 0$, $f_\lambda(x) = 1$ if $x \in \ell_{2,a}$, i.e. $f_\lambda(x) = \mathbf{1}_{x \in \ell_{2,a}} \rightarrow \lim_{\lambda \rightarrow 0} \int_{\mathbb{R}^\infty} f_\lambda d\mu = \int_{\ell_{2,a}} d\mu = m(\ell_{2,a}) = 1$.

i.e. we have a prob. measure μ that gives full measure to whole sep. space $\ell_{2,a}$ (isomorphic to L_2).

- It can be shown that in ℓ^2 w/ Gaussian measure $\mu = N(0, S_x = (a_{kn})_{k,n=1}^{\infty})$.

$$x \in \ell^2 \quad \text{if } \sum_{k=1}^{\infty} a_k < \infty.$$

i.e. infinite dim. sep. $x \in \ell^2$ has finite Gaussian measure

For Personal Use Only -bkwk2

Full covariance operators in infinite dimensions

matrix wrt infinite dim.

- In a Hilbert space H , the Gaussian measure μ w/ mean M and covariance operator C
 $\mu = N(M, C)$ has the rep. of linear operator C being trace class (nuclear), i.e. $\text{trace}(C) < \infty$.
- The trace operator in Hilbert space w/ on orthonormal basis is defined by

$$\boxed{\text{trace}(C) = \sum_{n=1}^{\infty} \langle c_n, C c_n \rangle}$$

Note for zero mean functions f, g on domain D s.t. $f, g \in H$, we define for $u \in U$,

$$\begin{aligned} \langle f, g \rangle &= E[\langle f, u \rangle \langle u, g \rangle] = E\left[\int_D \int_D f(x) u(x) u(y) g(y) dy dx\right] \\ &= E\left[\int_D f(x) \left(\int_D u(x) u(y) g(y) dy\right) dx\right] = \int_D f(x) \left[\int_D E[u(x) u(y)] g(y) dy \right] dx \\ &= \int_D f(x) \left(\int_D C(x, y) g(y) dy \right) dx. \end{aligned}$$

so we have $|c_{xy} = \int_D C(x, y) g(y) dy|$, $\langle c_n, c_n \rangle = \int_D c_n(x) \int_D C(x, y) g(y) dy dx$

+ $C(x, y)$ is the covariance function defining a Gaussian process.

EQUINODIMENSIONALITY OF GAUSSIAN MEASURES

- In finite dim, two Gaussian measures are absolutely continuous (equivalent) wrt each other.
e.g. In \mathbb{R} , $\mu = N(0, \sigma^2)$, $\mu_M = N(M, \sigma^2)$, both have the densities wrt Lebesgue measure
 \rightarrow for $\sigma \neq 0$, any $A \in \mathbb{R}$, $\mu(A) = 0 \Leftrightarrow \mu_M(A) = 0 \Leftrightarrow M(A) = 0$.
- However, in ℓ^2 , $\mu = N(0, C)$ and $\mu_V = N(V, C)$ are either equivalent or singular
wrt each other, depending on V . (Gaussian measure is quasi-invariant)
- The subspace containing all translations v for which the measures are absolutely continuous wrt each other is the image of the space ℓ^2 under the operator $C^{1/2}$.
(this space is the Cameron-Martin space / RKHS).

BAYES' RULE IN HILBERT SPACE

- In Hilbert space, the Radon-Nikodym derivative b/w a posterior measure μ^y and a prior measure μ^0 can define the corresponding likelihood function

$$\boxed{\frac{d\mu^y}{d\mu^0}(x) = L(y|x)}$$

- The reference measure for a Hilbert space is the Gaussian measure, so the prior measure over functions is $\mu^0 = N(M, C)$. (formal defn of Gaussian Process prior)

High dimensional MCMC

Gaussian random walk Metropolis Hastings (GRW-MH) algorithm

- For GRW-MH, the candidate sample v is generated by adding Gaussian noise $\xi \sim N(0, C)$ to the current sample u

$$v = u + \beta \xi$$

where β is the step size parameter

- The corresponding proposal function q is symmetric

$$q(u, v) = N(v | u, \beta^2 C) = N(u | v, \beta^2 C) = q(v, u)$$

The acceptance probability α therefore can be simplified to

$$\alpha(u, v) = \min\left\{\frac{\pi(v) q(v, u)}{\pi(u) q(u, v)}, 1\right\} = \min\left\{\frac{\pi(v)}{\pi(u)}, 1\right\}$$

- In an infinite dim Hilbert space H , the proposal is well-defined and random draws will have finite norm and converge if C is trace class

There is no Lebesgue measure in infinite dim \rightarrow replace $\pi(x)$ w/ $\pi^*(x)$.

- For the case $\mu^* = N(0, C)$ and $\frac{d\pi^*}{d\pi}(x) \propto \|y(x)\|^{-1} \exp(-\phi(x))$, the acceptance probability α becomes

$$\alpha(u, v) = \min\left\{\frac{\pi^*(v)}{\pi^*(u)}, 1\right\} = \min\left\{\frac{\|y(v)\| \mu^*(v)}{\|y(u)\| \mu^*(u)}, 1\right\} = \min\left\{\frac{\exp(\phi(v) - \frac{1}{2}\|C^{-1/2}v\|^2)}{\exp(\phi(u) - \frac{1}{2}\|C^{-1/2}u\|^2)}, 1\right\} = \min\left\{\exp(J(v) - J(u)), 1\right\},$$

where $J(x) = -\frac{1}{2}\langle x, Cx \rangle - \phi(x) = -\frac{1}{2}\|C^{1/2}x\|^2 - \phi(x)$. Note $\langle x, Cx \rangle = \|x\|_C^2$

- Let the eigenvalues of C be λ_i^2 . $u \sim N(0, C) \rightarrow u_i \sim N(0, \lambda_i^2) \rightarrow \frac{1}{\lambda_i} u_i \sim N(0, 1)$.

The eigenvalues of $C^{1/2}$ are $\frac{1}{\lambda_i}$, so $E[\|C^{1/2}u\|^2] = E\left[\sum_i \left(\frac{1}{\lambda_i} u_i\right)^2\right] = \sum_i \lambda_i^{-2} = \infty$.

\rightarrow GRW-MH does not work for infinite dim.

Preconditioned Crank-Nicholson (PCN) algorithm

- For PCN, the candidate sample v is a linear combination of the current sample u and Gaussian noise $\xi \sim N(0, C)$

$$v = \sqrt{1-\beta^2} u + \beta \xi$$

where β is the step size parameter

- The proposal function is given by $q(x, y) = N(y | \sqrt{1-\beta^2}x, \beta^2 C) = \exp\left(-\frac{1}{2\beta^2} \langle y - \sqrt{1-\beta^2}x, C^{-1}(y - \sqrt{1-\beta^2}x) \rangle\right)$, so
- $$-2\beta^2 \log \frac{q(v, u)}{q(u, v)} = \langle u, C^{-1}u \rangle - \langle v, C^{-1}v \rangle - 2\sqrt{1-\beta^2} \langle v, C^{-1}u \rangle + 2\sqrt{1-\beta^2} \langle v, C^{-1}v \rangle - (\frac{1}{2}\beta^2) \langle u, C^{-1}u \rangle + (\frac{1}{2}\beta^2) \langle v, C^{-1}v \rangle$$
- $$\therefore \frac{q(v, u)}{q(u, v)} = \exp\left(-\frac{1}{2} \langle u, C^{-1}u \rangle + \frac{1}{2} \langle v, C^{-1}v \rangle\right)$$

The acceptance probability α therefore can be simplified to

$$\alpha(u, v) = \min\left\{\frac{\pi^*(v) q(v, u)}{\pi^*(u) q(u, v)}, 1\right\} = \min\left\{\frac{\|y(v)\| \mu^*(v) \exp(-\frac{1}{2} \langle u, C^{-1}u \rangle)}{\|y(u)\| \mu^*(u) \exp(-\frac{1}{2} \langle v, C^{-1}v \rangle)}, 1\right\}$$

- For the case $\mu^* = N(0, C)$ and $\frac{d\pi^*}{d\pi}(x) \propto \|y(x)\|^{-1} \exp(-\phi(x))$, the acceptance probability α becomes

$$\alpha(u, v) = \min\left\{\frac{\|y(v)\|}{\|y(u)\|}, 1\right\} = \min\{\exp(d(u) - d(v)), 1\}$$

\rightarrow well defined for infinite dim, provided ϕ is well defined.

For Personal Use Only -bkwk2

Derivation of pCN algorithm.

- consider the stochastic DE w/ $\mu = N(0, C)$ as an invariant measure

$$du = -uds + \sqrt{2C} db$$

where b is standard brownian motion.

- applying the Trotter discretisation scheme yields a discrete-time Markov chain

$$v-u = -\delta((1-\theta)u + \theta v) + \sqrt{2\delta}\xi_0, \quad \theta \in [0,1]$$

where u is the current position, v is the next position, δ is the discrete time diff., and $\xi_0 \sim N(0, 1)$

Rearranging for v :

$$v = (1+\theta\delta)^{-1}(u - \delta(1-\theta)u + \sqrt{2\delta}\xi_0)$$

Setting $\theta = 1/2$:

$$v = (1 + \frac{\delta}{2})^{-1}(\frac{\delta + 1}{2}u + \sqrt{2\delta}\xi_0)$$

Let $\beta = \sqrt{2\delta}(1 + \frac{\delta}{2})^{-1}$:

$$v = \sqrt{\frac{\delta}{1-\delta}}u + \beta\xi_0$$

- It turns out the two measures defining the acceptance prob. are equivalent only for $\theta = \frac{1}{2}$

Taking the i -th coordinate of $v = (1+\theta\delta)^{-1}(u - \delta(1-\theta)u + \sqrt{2\delta}\xi_0)$, we have

$$v_i = \frac{1 - \delta(1-\theta)}{1 + \theta\delta} u_i + \frac{\sqrt{2\delta\lambda_i^2}}{1 + \theta\delta} g_i$$

where λ_i^2 are the eigenvalues of C and g_i is the standard Normal

The correspondingly expectation and variance are given by

$$\mathbb{E}[v_i] = 0 \quad \text{Var}[v_i] = \left(\frac{1 - \delta(1-\theta)}{1 + \theta\delta}\right)^2 \lambda_i^2 + \frac{2\delta\lambda_i^2}{(1 + \theta\delta)^2}$$

For two product Gaussian measures to be equivalent, we have the convergence criterion

$$\sum_{i=1}^{\infty} \left(\frac{\text{Var}[v_i]}{\text{Var}[u_i]} - 1 \right)^2 < \infty, \quad \text{r.e. } \frac{\text{Var}[v_i]}{\text{Var}[u_i]} = 1$$

$$\rightarrow \frac{\text{Var}[v_i]}{\text{Var}[u_i]} = \left(\frac{1 - \delta(1-\theta)}{1 + \theta\delta} \right)^2 + \frac{2\delta}{(1 + \theta\delta)^2} = 1 \rightarrow \theta = \frac{1}{2},$$

Langevin MCMC

Log concave and m -strongly log-concave

- A probability measure w/ density π is log-concave if

$$\log \pi(\lambda x + (1-\lambda)y) \geq \lambda \log \pi(x) + (1-\lambda) \log \pi(y), \quad \lambda \in [0,1]$$

More compactly, $\pi(\lambda x + (1-\lambda)y) \geq \pi(x)^\lambda \pi(y)^{1-\lambda}$

- A probability measure w/ density π is m -strongly log-concave if

$$\log \pi(\lambda x + (1-\lambda)y) \geq \lambda \log \pi(x) + (1-\lambda) \log \pi(y) + \frac{m\lambda(1-\lambda)}{2} \|x-y\|^2, \quad \lambda \in (0,1)$$

This is a stronger notion that req. there not be flat regions and the max. is unique

e.g. Gaussians - log-concave = quadratic function

For Personal Use Only -bkwk2

The Langevin SDE

- The Langevin SDE is given by

$$dx_t = -\nabla U(x_t) dt + \sqrt{2} dB_t$$

where $U: \mathbb{R}^d \rightarrow \mathbb{R}$ is a potential and $\{B_t\}_{t \geq 0}$ is Brownian motion

- The Langevin SDE has the stationary measure

$$\pi \propto \exp(-U(x))$$

i.e. to bound the π , we set $U(x) = -\log \pi(x) \rightarrow dU(x) = -\nabla \log \pi(x)$, so we have

$$dx_t = \nabla \log \pi(x_t) dt + \sqrt{2} dB_t$$

- Before reaching the stationary measure, $x_{t \rightarrow \infty}$ satisfies the Fokker-Planck equation

$$\frac{\partial p_t}{\partial t} = -\nabla \cdot (p_t \nabla \log \pi) + \nabla^2 p_t$$

Asymptotically, we have $\partial p_t = 0$, (as $t \rightarrow \infty$).

$$\nabla^2 p_t = \nabla \cdot (p_t \nabla \log \pi) \rightarrow \nabla p_t = p_t \nabla \log \pi \rightarrow \frac{1}{p_t} \nabla p_t = \nabla \log \pi \rightarrow \nabla \log p_t = \nabla \log \pi.$$

$$\lim_{t \rightarrow \infty} p_t = \pi$$

→ Simulating the SDE provides samples from π in the limit.

- the Langevin SDE has many nice properties:

↳ For a m -strongly log-concave π , the convergence rate to the stationary measure π is exponentially fast,

$$W_2(\pi_t, \pi) \leq e^{-\gamma t} \left\{ \|x - \bar{x}\| + (\gamma m)^{1/2} \right\} \quad \text{for } \pi_t = \delta_{\bar{x}}, \bar{x} = \arg \max_{x \in \mathbb{R}^d} \log \pi(x)$$

where W_2 is a distance on the space of prob. measures.

Metropolis-adjusted Langevin algorithm (MALA)

- The continuous SDE can't be simulated exactly → use discretization schemes, but it introduces bias → we correct for the bias via a Metropolis accept/reject step.
- consider first-order Euler discretisation, using $x_0 \sim \pi_0$, we have

$$x_{k+1} = x_k + \gamma \nabla \log \pi(x_k) + \sqrt{2\gamma} w_{k+1}$$

where $\{w_k\}_{k \geq 0}$ are standard Normal and $\gamma > 0$ is the step size.

- Given the state of the Markov chain x_k , propose candidate sample \hat{x}_{k+1} ,

$$\hat{x}_{k+1} = x_k + \gamma \nabla \log \pi(x_k) + \sqrt{2\gamma} w_{k+1}$$

and accept w/ probability

$$\alpha(x_k, \hat{x}_{k+1}) = \min \left\{ \frac{\pi(\hat{x}_{k+1}) q(x_{k+1}, x_k)}{\pi(x_k) q(x_k, \hat{x}_{k+1})}, 1 \right\}$$

where $q(x, x') = N(x' | x + \gamma \nabla \log \pi(x), 2\gamma) \propto \exp(-\frac{1}{4\gamma} \|x' - x - \gamma \nabla \log \pi(x)\|^2)$.

- The step size γ has to be tuned s.t. $\alpha = 0.574$.

- MALA may suffer in high dimensions d – acceptance prob. α may get exponentially smaller in d , which results in poor mixing.

For Personal Use Only -bkwk2

Unadjusted Langevin algorithm (ULA)

- Run a Langevin SDE discretisation w/o the Metropolis step, i.e.

$$X_{k+1} = X_k + \gamma \nabla \log \pi(X_k) + \sqrt{2\gamma} W_{k+1}$$

The chain converges to a stationary measure π_T which is not the target measure π , but performs well in high dim. (Metropolis step avoided)

- consider the example of a Gaussian target $\pi(x) = N(x|\mu, \Sigma)$.

The Langevin SDE in this case is the Ornstein-Uhlenbeck (OU) process,

$$dX_t = -\Sigma^{-1}(X_t - \mu) dt + \sqrt{2} dB_t$$

The ULA scheme is given by

$$X_{k+1} = X_k - \gamma \Sigma^{-1}(X_k - \mu) + \sqrt{2\gamma} W_{k+1}$$

Under appropriate conditions, the chain has the invariant measure

$$\pi_T = N(\mu, \Sigma(I - \frac{\gamma}{2}\Sigma^{-1})^{-1}) \quad \leftarrow \begin{array}{l} \text{proof in EP, } X_{k+1} \text{ is Gaussian} \\ \text{consider } E[X_{k+1}], \text{Var}[X_{k+1}] \end{array}$$

→ simulates the correct mean but introduces bias to the variance

- For m-strongly log-concave π , the convergence rate of the ULA w/ respect to π_0

$$W_2(\pi_0, \pi) \leq (1-\gamma)^k W_2(\pi_0, \pi) + 1.65(4m)(\gamma k)^{1/2}, \quad \gamma \leq 2/mL$$

where we have L -Lipschitz gradients and W_2 is a distance on the space of prob. measures

* the asymptotic bias is of order $O(\gamma^{1/2})$ - can be made arbitrarily small w/ small γ .

Langevin schemes for Bayesian inference

- consider target distributions of the form $\pi(x) \propto p(x) \prod_{i=1}^n L(y_i|x)$, where $p(x)$ is the prior over x and $L(y_i|x)$ is the likelihood for data vector y_i . usually $n \gg 1$.

- If we apply MALA, we run into two problems:

↳ Gradient computation for the proposal is too costly. $O(n)$ → bottleneck.

$$\nabla \log \pi(x) = \nabla \log p(x) + \sum_{i=1}^n \nabla \log L(y_i|x)$$

↳ Acceptance prob. computation too costly. $O(n)$ → avoided w/ Metropolis-free schemes.

$$\alpha = \min \left\{ \frac{\pi(X_{k+1}) q(X_{k+1}, X_k)}{\pi(X_k) q(X_k, X_{k+1})}, 1 \right\}, \quad \pi(x) \propto p(x) \prod_{i=1}^n L(y_i|x)$$

- We can estimate the gradient cheaply if we subsample data - random minibatch.

$$\widehat{\nabla \log \pi(x_k)} = \nabla \log p(x_k) + \frac{1}{|I_k|} \sum_{i \in I_k} \nabla \log L(y_i|x_k)$$

where $I_k \subset \{1, \dots, n\}$ is the index set of iteration k , and $|I_k|$ is the size of the mini batch,

- If we run ULA w/ $\widehat{\nabla \log \pi(x_k)}$, we have stochastic gradient Langevin dynamics (SGLD)

$$X_{k+1} = X_k + \gamma \widehat{\nabla \log \pi(x_k)} + \sqrt{2\gamma} W_{k+1}$$

where $E[\widehat{\nabla \log \pi(x_k)}] = \nabla \log \pi(x_k)$ and each iteration costs $O(|I_k|)$

- For m-strongly log-concave π , we have similar convergence guarantees of ULA

$$W_2(\pi_0, \pi) \leq O(e^{-\gamma M^2} + \gamma \frac{1}{2})$$

* we just need to ensure the variance of the gradient estimates are controlled.