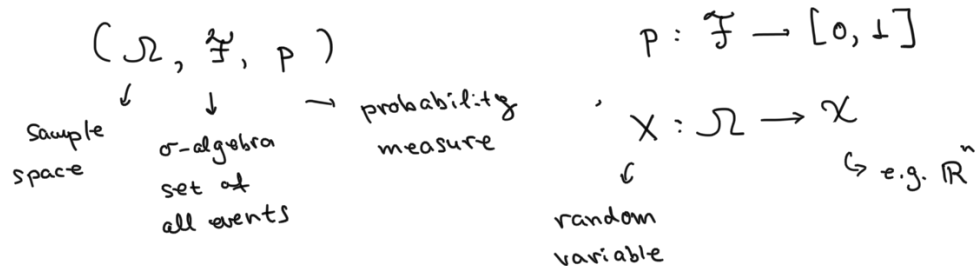


Lecture #1 : Primer on probability and statistics

Notation :

- Lower case : e.g. x, y, w, z $\left\{ \begin{array}{l} \text{event realizations, data points,} \\ \text{column vectors (row vectors e.g. } w^T) \\ \text{functions, probability measures, etc.} \end{array} \right.$
- Upper case : e.g. X, Y, A, B $\left\{ \begin{array}{l} \text{random variables} \\ \text{matrices} \\ \text{some functions (e.g. cdf } F) \end{array} \right.$
- Caligraphics : $\{ \text{sets, operators, e.g. } \mathcal{F}, \mathbb{E} \}$

Probability space :



Discrete random variables :

A discrete r.v. X is an event taking values in a finite or countably infinite discrete space \mathcal{X} .

We will denote the probability of $X=x$ as $p(X=x)$

or $\boxed{p(x)}$ (more formally : $p(\{\omega \in \Omega : X(\omega) = x\})$)

- $\left\{ \begin{array}{l} \text{i) } 0 \leq p(x) \leq 1 \\ \text{ii) } \sum_{x \in \mathcal{X}} p(x) = 1 \end{array} \right.$

, here $p(x)$ is called the mass function of X .

$$(iii) \quad p(\{\emptyset\}) = 0$$

probability
(pmf)

Continuous random variables: (events that take continuous values)

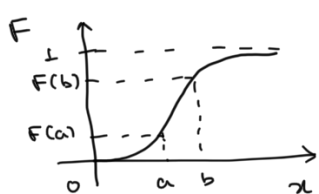
Define the events: $A = (X \leq a)$, $B = (X \leq b)$, $W = (a < X \leq b)$

Observe, $B = A \cup W$, $A \cap W = \{\emptyset\}$

$$P(B) = P(A) + P(W) - P(A \cap W)$$

$$\Rightarrow P(W) = P(B) - P(A)$$

Define a function $F(x) := P(X \leq x) \rightarrow$ cumulative distribution function of (cdf)



$$P(W) = P(a < X \leq b) = F(b) - F(a)$$

properties of a cdf

$$(i) \quad 0 \leq F(x) \leq 1$$

$$(ii) \quad \lim_{x \rightarrow -\infty} F(x) = 0, \quad \lim_{x \rightarrow +\infty} F(x) = 1$$

$$(iii) \quad x \leq y \Rightarrow F(x) \leq F(y) \quad \text{monotonically non decreasing function}$$

If the cdf is differentiable, define $p(x) := \frac{d}{dx} F(x)$

\hookrightarrow probability density function (pdf) of X

By definition:

$$P(W) = P(a < X \leq b) = \int_a^b p(x) dx$$

Properties of a pdf:

$$(i) \quad p(x) \geq 0$$

$$(ii) \quad \int_{-\infty}^{\infty} p(x) dx = 1$$

$$(iii) \quad \int p(x) dx = P(X \in A)$$

(*) We require that

$p(x) \geq 0$, it is possible

that $p(x) \geq 1 \forall x$, as

long as long as:

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

$$x \in A$$

Lecture #2

Quantiles :

$$F(x) := P(X \leq x)$$

If the cdf of a continuous random variable X is a monotonically non-decreasing function, then it has an inverse. Then $F^{-1}(\alpha)$ is the value x_α such that $P(X \leq \underline{x_\alpha}) = \alpha$.

This probability is called the α -quantile of X .

e.g. $F^{-1}(0.5)$ is called the median of the distribution.

In general, we can use the inverse cdf to compute tail area probabilities.

Mean / Expected value / 1st-order moment :

• Discrete case : $\mu = E[X] := \sum_{x \sim p(x)} x \cdot p(x)$

• Continuous case : $\mu = E[X] = \int_{\mathcal{X}} x p(x) dx$

Properties of expectation:

• $E[c] = c$

• $E[c f(x)] = c E[f(x)] = c \int_{\mathcal{X}} f(x) p(x) dx$

• $E[f(x) + g(x)] = E[f(x)] + E[g(x)]$

Variance / 2nd-order moment :

$$\sigma^2 = \text{Var}[X] := \mathbb{E}[(X-\mu)^2] = \int_{\mathcal{X}} (x-\mu)^2 p(x) dx$$

Remark: $\boxed{\mathbb{E}[X^2] = \mu^2 + \sigma^2}$

$$\begin{aligned} \sigma^2 &:= \int_{\mathcal{X}} (x-\mu)^2 p(x) dx = \underbrace{\int_{\mathcal{X}} x^2 p(x) dx}_{\mathbb{E}[X^2]} + \mu^2 \int_{\mathcal{X}} p(x) dx - 2\mu \underbrace{\int_{\mathcal{X}} x p(x) dx}_{\mu} \\ &= \mathbb{E}[X^2] - \mu^2 \end{aligned}$$

Standard deviation:

$$\sigma = \text{Std}[X] = \sqrt{\text{Var}[X]} = \sqrt{\sigma^2}$$

• Degenerate pdf: $\lim_{\sigma \rightarrow 0} \mathcal{N}(\mu, \sigma^2) = \delta(x-\mu),$

$$\delta(x) = \begin{cases} \infty, & \text{if } x=0 \\ 0, & \text{if } x \neq 0 \end{cases}$$

↓
Dirac delta

$$\int_{-\infty}^{+\infty} \delta(x) dx = 1$$

• Shifting property: $\int_{-\infty}^{+\infty} \underbrace{\delta(x-\mu)}_{\substack{\neq 0 \\ 0, \text{ only when } x=\mu}} p(x) dx = f(\mu)$

• Empirical measure / distribution:

Given some observations $\mathcal{D} := \{x_1, \dots, x_n\}$ we define:

$$p_e(\mathcal{D}) := \frac{1}{n} \sum_{i=1}^n \delta_{x_i}(\mathcal{D}), \quad \delta_{x_i}(\mathcal{D}) \text{ is the Dirac measure}$$

$$\delta_{x_i}(\mathcal{D}) = \begin{cases} 1, & x_i \in \mathcal{D} \\ 0, & x_i \notin \mathcal{D} \end{cases}$$

NOTE: ... also associate weights \pm .

we can also assign weights to each observation:

$$p_{\theta}(\mathcal{D}) = \sum_{i=1}^n w_i \delta_{x_i}(\mathcal{D}) \quad , \quad 0 \leq w_i \leq 1 \quad , \quad \sum_{i=1}^n w_i = 1$$

Joint distributions: $p(x_1, x_2, \dots, x_d)$, $x \in \mathbb{R}^d$
 $x = (x_1, x_2, \dots, x_d)$

Covariance:

The covariance between two random variables X, Y measures the degree to which X and Y are linearly related.

$$\begin{aligned} \text{cov}[X, Y] &:= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])^T] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \end{aligned}$$

If $x \in \mathbb{R}^d$ is a d -dimensional random vector, then its covariance matrix is a symmetric and positive definite matrix:

$$\text{cov}[x] = \text{cov}[x, x] = \mathbb{E}[(x - \mathbb{E}[x])(x - \mathbb{E}[x])^T]$$

$$\begin{bmatrix} \text{Var}[x_1] & \text{cov}[x_1, x_2] & \dots & \text{cov}[x_1, x_d] \\ \vdots & \text{Var}[x_2] & & \vdots \\ \vdots & & \ddots & \vdots \\ \text{cov}[x_d, x_1] & \dots & \dots & \text{Var}[x_d] \end{bmatrix}_{d \times d} = \overset{\substack{\text{covariance} \\ \text{matrix}}}{\Sigma} = \overset{-1}{\Lambda} \quad \hookrightarrow \text{precision matrix}$$

- Covariances can take values between zero and infinity.

Sometimes it is more preferable to work with a normalized

measure with a finite upper bound:

$$\text{Pearson correlation : } \text{corr}[X, Y] := \frac{\text{cov}[X, Y]}{\sqrt{\text{Var}[X] \text{Var}[Y]}}$$

$$-1 \leq \text{corr}[X, Y] \leq 1.$$

Specifically one can show that $\text{corr}[X, Y] = 1$ iff

$$Y = \alpha X + b, \text{ for some } \alpha, b.$$

Independence : X, Y are (unconditionally) independent if

$$X \perp Y \iff P(X, Y) = P(X)P(Y)$$

X, Y are conditionally independent if

$$X \perp Y | Z \iff P(X, Y | Z) = P(X | Z)P(Y | Z)$$

If X, Y are independent : $\text{cov}[X, Y] = 0$ hence X, Y are also uncorrelated.

Caution : The opposite is not always true!

$$\text{Entropy : } H[X] := - \int_{x \sim P(X)} P(x) \log P(x) dx$$

$$\text{Relative entropy : } x \sim \underline{P(x)}, y \sim \underline{Q(y)} \quad \text{KL}[Q \| P]$$

$$H[X|Y] = - \int \log \frac{P(x)}{Q(y)} P(x) dx = \text{KL}[P \| Q]$$

↓
Kullback-Leibler divergence

Mutual information :

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) = H(X, Y) - H(X) - H(Y)$$

$$I(X, Y) = \iint_{x, y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} dx dy \quad (= 0 \text{ if } X, Y \text{ are independent})$$