

Probability Review :

A random variable, e.g. X , is a variable whose possible values are numerical outcomes of a random phenomenon.

Discrete rvs: takes countable number of distinct values, e.g. 2, 3, 4, ... Ex: Flipping a coin $H \rightarrow X=1$
 $T \rightarrow X=0$
 probability distribution: list of probabilities assigned to each possible value. $P(X=1) = P(H) = p_1$
 $P(X=0) = P(T) = p_2$
 $\text{probability distribution} = \text{probability mass function}$

Continuous rvs: takes infinite number of possible values, e.g., temperature of Boston $X \in (-\infty, +\infty)$
 cont. rv is not defined at specific values, it is defined over an interval of values and represented by area under curve. $P(X=x) = 0$, $P(x_1 \leq X \leq x_2) = \int_{x_1}^{x_2} p(x) dx$
 $\text{density function} \rightarrow p(x) \geq 0$
 $\int p(x) dx = 1$

Joint distribution: For collection of random variables e.g. X_1, Y or X_1, X_2, \dots, X_N
 $p(X, Y)$ or $p(X_1, X_2, \dots, X_N)$
 Boston temperature Flipping a coin Flipping a coin N times

Conditioning: $P(X=x | Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)}$, $\sum_x P(X=x | Y=y) = 1$
 but not necessarily $\sum_y P(X=x | Y=y) \neq 1$

Chain rule of probability: $P(X, Y) = \frac{P(X, Y)}{P(Y)} \rightarrow P(X, Y) = P(X|Y)P(Y)$ or $P(Y|X)P(X)$

$$P(X_1, X_2, \dots, X_N) = P(X_1) P(X_2 | X_1) P(X_3 | X_1, X_2) \dots P(X_N | X_1, \dots, X_{N-1})$$

Marginalization: Given joint pd. $P(X, Y) \rightarrow P(X) = \sum_y P(X, Y=y) = \sum_y P(Y=y | X) P(X)$
 $= P(X) \times 1 = P(X)$

$$P(X, Y, Z) \rightarrow P(Y) = \sum_x \sum_z P(X=x, Y, Z=z)$$

Bayes Rule: $P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$ e.g. $P(D|D) = \frac{P(D|D)P(D)}{P(D)}$ MAP

Independence: X and Y are independent if observing X does not affect information about Y

$$X \perp\!\!\!\perp Y \Leftrightarrow P(Y|X) = P(Y) \equiv P(X, Y) = P(X)P(Y) \quad \text{e.g. } X = H, T$$

$Y = \text{Boston temp.}$

Conditional Independence : X and Y are conditionally independent given Z iff $P(X, Y | Z) = P(X | Z) P(Y | Z)$
 $X \perp\!\!\!\perp Y | Z$

e.g. Height and Vocabulary are not independent (small child has a limited vocab.)

but given age they become independent : Height $\perp\!\!\!\perp$ Vocabulary $|$ Age

Expectation: Given X being a rv, any function $f(x)$ is also a rv.

$$E[f(X)] = \sum_x f(x) p(x)$$

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

$$E[f(X) | Y=y] = \int f(x) p(x|y) dx \quad \text{or} \quad \sum_x p(x|y) f(x)$$

IID: independent and identically distributed rvs are mutually independent and all come from the same distribution: E.g. flip a coin N times.

Probability Model: probability distribution is used as models of observed data

Ex) We flip a coin N times and record X_1, X_2, \dots, X_N

We measure Boston city temperature on Jan 1, each year X_1, \dots, X_N

These prob. dist. are parameterized and the goal is to learn the parameters of the p.d.

Ex) Flip a coin N times : X_1, X_2, \dots, X_N : $P(X_i=1) = \theta$, $P(X_i=0) = 1-\theta \Rightarrow P(x_i) = \theta^{x_i} (1-\theta)^{1-x_i}$

Given Dataset : $P_\theta(x_1, x_2, \dots, x_N) \stackrel{\text{iid}}{=} \prod_{i=1}^N P_\theta(x_i)$

$$x_1 = x_1, \dots, x_N = x_N$$

$$= \prod_{i=1}^N \theta^{x_i} (1-\theta)^{1-x_i} = \theta^{\sum_{i=1}^N x_i} (1-\theta)^{N - \sum_{i=1}^N x_i} = \text{likelihood function}$$

* Given x_1, \dots, x_N , determine parameters of the prob. model on data.

How, determine θ ? $\max_{\theta} P_\theta(x_1, \dots, x_N)$ for what value of θ , the observations are most likely?

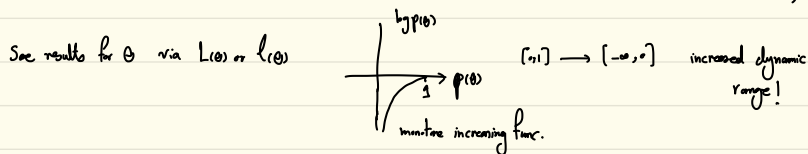
e.g. $X_1=1, X_2=1, \dots, X_N=1 \rightarrow$ if $\theta=0 \rightarrow P(x_i=1) = 0 \rightarrow$ we never observe 1!
if $\theta=1 \rightarrow P(x_i=1) = 1 \rightarrow$ we always observe 1!

To determine parameters of the model, we use maximum likelihood:

$$L(\theta) \triangleq p_{\theta}(x_1, \dots, x_N) \rightarrow \max_{\theta} L(\theta) = \max_{\theta} p_{\theta}(x_1, \dots, x_N)$$

Generally, it is more convenient to take \log (simplifies operations, increases dynamic range)

$$l(\theta) \triangleq \log p_{\theta}(x_1, \dots, x_N) \rightarrow \max_{\theta} l(\theta) = \max_{\theta} \log p_{\theta}(x_1, \dots, x_N) \rightarrow \text{maximum log-likelihood}$$



$$l(\theta) = \sum_{i=1}^N x_i \log \theta + (N - \sum_{i=1}^N x_i) \log (1 - \theta) \rightarrow \text{is it concave?}$$

$$\frac{\partial l}{\partial \theta} = \sum_{i=1}^N x_i / \theta - \frac{N - \sum_{i=1}^N x_i}{1 - \theta} = 0 \rightarrow \hat{\theta}_{ML} = \frac{\sum_{i=1}^N x_i}{N} = \text{sample mean}$$

Gaussian MLE: $x_i \sim \mathcal{N}(x_i; \mu, \sigma^2) \Rightarrow p(x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$

$$\theta = \{\mu, \sigma^2\}: p_{\theta}(x_1, \dots, x_N) \stackrel{\text{iid}}{=} \prod_{i=1}^N p_{\theta}(x_i) = \frac{1}{(2\pi\sigma^2)^{N/2}} e^{-\frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}}$$

$$l(\theta) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^2}$$

$$\frac{\partial l}{\partial \theta} = \begin{pmatrix} \frac{\partial l}{\partial \mu} \\ \frac{\partial l}{\partial \sigma^2} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \rightarrow \frac{\partial l}{\partial \mu} = \sum_{i=1}^N (\mu - x_i) / \sigma^2 = 0 \rightarrow \hat{\mu}_{ML} = \frac{\sum_{i=1}^N x_i}{N}$$

$$\frac{\partial l}{\partial \sigma^2} = -\frac{N}{2} \times \frac{1}{\sigma^2} + \frac{\sum_{i=1}^N (x_i - \mu)^2}{2\sigma^4} = 0 \rightarrow \hat{\sigma}_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}_{ML})^2$$

empirical mean and variance $\rightarrow \hat{\mu}_{ML}, \hat{\sigma}_{ML}^2$