

1 Basic definitions

Definition 1.1. Support The set of values a random variable can take (set of all the numerical realizations of outcomes) e.g. for a binary random variable X , it's support is $R_X = \{0, 1\}$

Definition 1.2. Sample Point In an experiment, a sample point is one of the possible outcomes of experiment denoted by ω

Definition 1.3. Independent Events Two events A and B are independent if and only if

$$P(A \cap B) = P(A) * P(B) \quad (1)$$

Knowing that B has occurred, does not make us change our belief of the probabilities of A , and vice versa.

Definition 1.4. Jointly/Mutually Independent Events Let E_1, \dots, E_n be n events.

E_1, \dots, E_n are jointly independent (or mutually independent) if and only if for any sub-collection of k events ($k \leq n$) E_{i1}, \dots, E_{ik} :

$$P\left(\bigcap_{j=1}^k E_{ij}\right) = \prod_{j=1}^k P(E_{ij}) \quad (2)$$

Definition 1.5. Conditional Independence of Events Let there be three events: A, B, C . We say that B and C are independent, given the condition that event A has occurred. Before knowing that A has occurred, they might not have been independent.

Definition 1.6. Random Variable A random Variable X is a function from sample space Ω to set of real numbers \mathbb{R} , i.e $X : \Omega \rightarrow \mathbb{R}$.

The real number $X(\omega)$ associated with sample point $\omega \in \Omega$ is called a realization of the random variable. The set of all possible realizations is called the support and denoted by R_X

Definition 1.7. Independence of Random Variables Two random Variables X and Y are independent of each other, if the joint PMF of X and Y , satisfies following condition:

$$p_{X,Y} = p_X(x) \cdot p_Y(y) \quad \text{is true for } \forall x \in X, y \in Y \quad (3)$$

Understanding it in simple terms, for the 2-d table for PMF, any entry can be known by multiplying corresponding values in marginal probabilities $p_X(x)$ and $p_Y(y)$, if they are independent, and vice versa.

Definition 1.8. Probability Mass Function The PMF of a discrete random variable X is a function $p_X : \mathbb{R} \rightarrow [0, 1]$ such that

$$p_X = P(X = x) \quad \forall x \in \mathbb{R} \quad (4)$$

where $P(X = x)$ is the probability of realization of random variable X will be equal to x . Basically PMF is numerical realizations \rightarrow respective Probabilities.

Definition 1.9. Distribution Function/Cumulative Distribution Function If X is a random variable, its distribution/cdf is a function $F_X : \mathbb{R} \rightarrow [0, 1]$ such that

$$F_X(x) = P(X \leq x) \quad \forall x \in \mathbb{R} \quad (5)$$

Definition 1.10. Expected Value of a Random Variable The expected value of random variable X is the weighted average of the values that X can take on where each possible realization value is weighted by its respective probability i.e.

$$\mathbf{E}[X] = \sum_{x \in R_X} xp_X(x) \quad (6)$$

In other words, it is the weighted average of all possible numerical outcomes (in the support) weighted with respect to their respective probabilities

Definition 1.11. Expected value of function of Random Variable Expected value of function $g(Y)$ say $g(Y) = Y^2$:

$$\mathbf{E}[g(Y)] = \sum_{y \in R_Y} g(y)p(y) \quad (7)$$

Definition 1.12. Linearity of Expectation If X is a random variable and $a \in \mathbb{R}$ is a constant, then

$$\mathbf{E}[aX] = a\mathbf{E}[X] \quad (8)$$

In a more general setting:

If Y is a random variable such that $Y = a + bX$, where X is a random variable, and a and b are constants,

$$\begin{aligned} \mathbf{E}[Y] &= \mathbf{E}[a] + \mathbf{E}[bX] \\ &= a + b\mathbf{E}[X] \\ \mathbf{E}[Y] &= a + b\mathbf{E}[X] \end{aligned} \quad (9)$$

Definition 1.13. Conditional expectation of Random Variable Given random variables X and Y , The conditional expectation of X given $Y = y$ is the weighted average of values that X can take on, where each possible value is weighted by its respective conditional probability (conditional on the information $Y = y$), denoted by $E[X|Y = y]$

$$\mathbf{E}[X|Y = y] = \sum_{x \in R_x} xp_{X|Y=y}(x) \quad (10)$$

Another fact worth noting is that $\mathbf{E}[X|Y = y]$ when y is known is a number, but when y is unknown, $\mathbf{E}[X|Y]$ acts like a random variable that is a function of random variable Y i.e. $\mathbf{E}[X|Y] = g(Y)$

Definition 1.14. Total Expectation/Iterated Expectation The law of total expectation/law of iterated expectation/tower rule/smoothing theorem/adam's law among other names states that if, X is a random variable whose expected value $E[X]$ is defined, and Y is any random variable on same probability space then

$$\mathbf{E}[X] = \mathbf{E}[\mathbf{E}[X|Y]] \quad (11)$$

This comes from the fact that if there are A_i is a finite countable partition of the sample space, then

$$\mathbf{E}[X] = \sum_i \mathbf{E}[X|A_i]P[A_i] \quad (12)$$

In other words, the first equation is also summarized as, the expected value of a conditional expectation (an r.v.) is its unconditional expected value.

Definition 1.15. Geometric Random Variable Given that phenomena has a sequence of independent trials, there are two possible outcomes for each trial and the probability of success p is same for every trial, The geometric random variable X is the count of number of failures before the first success. i.e.

$Pr(X = k)$ means the first success happened on the k^{th} trial. The support $R_X = 1, 2, 3, \dots$. Don't confuse this with the outcomes of a single trial. We should really think that experiment is an aggregate process.

The **pmf of X** is looks like exponentially decreasing probability that adds up in total to 1.

$$P(X = k) = (1 - p)^k p \quad (13)$$

Definition 1.16. Quantile The p -Quantile of a Random Variable X is given $p \in (0, 1)$ and X having cumulative distribution function(cdf) to be $F_X(x)$, the p -Quantile is:

$$Q_X(p) = \min\{x \in R : F_X(x) \geq p\} \quad (14)$$

This translates roughly equivalently to the definition of percentiles. e.g 0.5-quantile is a median. 0.25-Quantiles are quartiles etc.

Definition 1.17. Deviation of a Random Variable The Deviation of a Random Variable is its difference from its mean value/Expected value. It is denoted by

$$\bar{X} = X - \mathbf{E}[X] \quad (15)$$

Deviation of Random Variable is also a random variable in it's own respect.

Definition 1.18. Variance It is a measure of dispersion of a random variable (expected squared deviation from expectation). Let X be a random variable. The variance of X , denoted by $Var[X]$ is defined as follows:

$$\mathbf{Var}[X] = \mathbf{E}[(X - \mathbf{E}[X])^2] = \mathbf{E}[\bar{X}^2] = \mathbf{E}[X^2] - \mathbf{E}[X]^2 \quad (16)$$

You can also think of it as mean of the square of the deviations of original random variable.

One would wonder why $\mathbf{E}[(X - \mathbf{E}[X])^2]$ is used instead of $\mathbf{E}[X - \mathbf{E}[X]]$, The reason is $X - \mathbf{E}[X]$ is signed distance, here is a short proof

$$\mathbf{E}[X - \mathbf{E}[X]] = \mathbf{E}[X - k] \quad (17)$$

$$= \mathbf{E}[X] - \mathbf{E}[k] \quad (18)$$

$$= k - k \quad (19)$$

$$= 0 \quad (20)$$

Definition 1.19. Standard Deviation It is a measure of dispersion of a random variable. Let X be a random variable. The standard deviation of X , denoted by $stdev[X]$ or $std[X]$ is defined as follows:

$$stdev[X] = \sqrt{Var[X]} \quad (21)$$

You can also think of it as rms deviation.

Definition 1.20. Covariance It is a measure of association between two random variables. Covariance of two random variables X and Y is , provided Expected values are well defined,

$$Cov[X, Y] = \mathbf{E}[(X - \mathbf{E}[X]) * (Y - \mathbf{E}[Y])] = \mathbf{E}[\bar{X} * \bar{Y}] \quad (22)$$

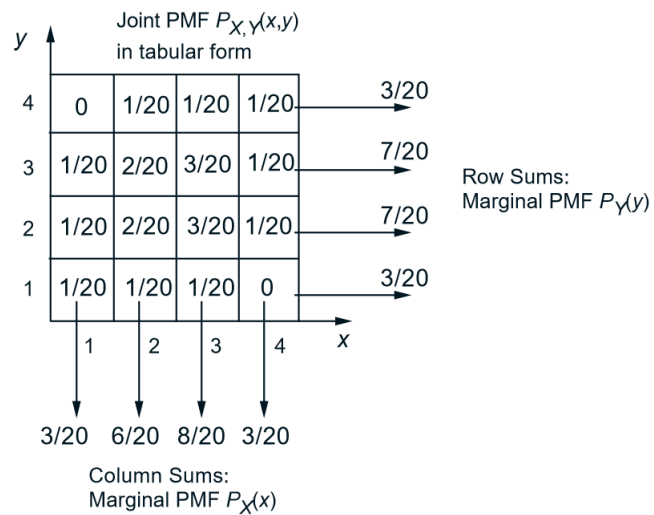
where we say deviation of X is $\bar{X} = X - \mathbf{E}[X]$ and

deviation of Y is $\bar{Y} = Y - \mathbf{E}[Y]$

or in other words, covariance is expectation of product of deviations.

Definition 1.21. Joint PMFs Joint probabilities tell how likely it is that certain X 's go together with certain Y 's. Notation is:

$$p_{X,Y}(x, y) = P(X = x \text{ and } Y = y) \quad (23)$$



It looks like following:

Useful in statistical studies that try to relate random variables to each other.

Some properties are:

- $\sum_x \sum_y p_{X,Y}(x, y) = 1$
- Finding PMF for single random variable is simple: $p_X(x) = \sum_y p_{X,Y}(x, y)$, these are also known as marginal PMFs.
- Conditional expectation means locking down a row/column as a new universe, i.e. $p_{X|Y}(x|y) = P(X = x|Y = y) = p_{X,Y}(x, y)/p_Y(y)$

Definition 1.22. Random Vector It is a multidimensional generalisation of the concept of Random Variable. Associated probability functions have the word "joint" in front of them e.g. Joint PMF, Joint cdf etc correspond to a Random Vector.

2 Logistic

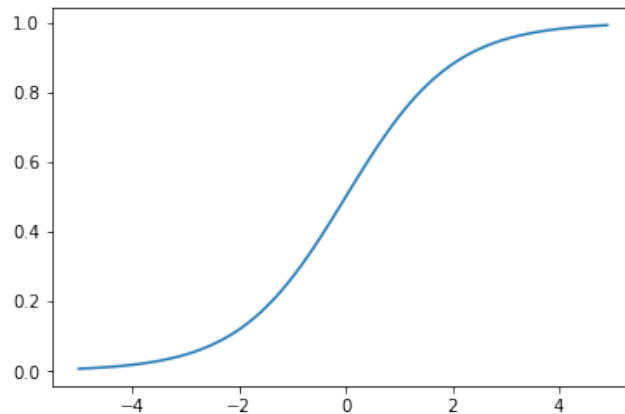
A logistic function or logistic curve is a common 'S' shaped curve with equation:

$$f(x) = \frac{L}{1 + e^{-k(x-x_0)}} \quad (24)$$

where

- e = natural log
- x_0 = x-value of sigmoid's midpoint
- L = the curve's maximum value
- k = the steepness of the curve

2.1 Standard Logistic



The standard logistic function is the logistic function with parameters given ($k = 1, x_0 = 0, L = 1$) i.e.

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (25)$$

which when plotted looks like

Why is logistic function so important? Because it can take any real input x , ($x \in R$), whereas the output always takes values between 0 and 1, and hence is interpretable as probability.

2.2 Derivative of logistic

Let's denote the sigmoid function as $\sigma(x) = \frac{1}{1 + e^{-x}}$.

The derivative of the sigmoid is $\frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x))$.

Here's a detailed derivation:

$$\frac{d}{dx}\sigma(x) = \frac{d}{dx} \left[\frac{1}{1 + e^{-x}} \right] \quad (26)$$

$$= \frac{d}{dx} (1 + e^{-x})^{-1} \quad (27)$$

$$= -(1 + e^{-x})^{-2}(-e^{-x}) \quad (28)$$

$$= \frac{e^{-x}}{(1 + e^{-x})^2} \quad (29)$$

$$= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} \quad (30)$$

$$= \frac{1}{1 + e^{-x}} \cdot \frac{(1 + e^{-x}) - 1}{1 + e^{-x}} \quad (31)$$

$$= \frac{1}{1 + e^{-x}} \cdot \left(\frac{1 + e^{-x}}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right) \quad (32)$$

$$= \frac{1}{1 + e^{-x}} \cdot \left(1 - \frac{1}{1 + e^{-x}} \right) \quad (33)$$

$$= \sigma(x) \cdot (1 - \sigma(x)) \quad (34)$$

2.3 Cost function of logistic regression

Prediction function is sigmoid of weighted sum of input features:

$$\hat{y} = \sigma(\theta^T \cdot x) \quad (35)$$

The log-loss function for a particular prediction is defined as $L(\hat{y}, y)$:

$$L(\hat{y}, y) = y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}) \quad (36)$$

Cost function for all examples is average of log-loss for each example i of the m examples:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})] \quad (37)$$