

Probability and Statistics

Omid Safarzadeh

January 30, 2022

Table of contents

- 1 Moments
 - Expected Value
 - Variance
- 2 Covariance and Correlation
 - Variance of Sums of Random Variables
- 3 Moment Generating Functions
 - Normal mgf
- 4 Matrix Notation for Moments
- 5 Distributions
 - Discrete Distribution
 - Discrete Uniform Distribution
 - Binomial Distribution
 - Poisson Distribution
 - Continuous Distribution
 - Uniform Distribution
 - Exponential Distribution
 - Normal Distribution
 - Lognormal Distribution
 - Laplace distribution
- 6 Reference

***Acknowledgement:** This slide is prepared based on Casella and Berger, 2002

Definition 1.1

For each of integer n , the n^{th} moment of X is

$$\mu'_n = E[X^n].$$

The n^{th} **central moment** of X , μ_n , is

$$\mu_n = E[(X - \mu)^n],$$

where $\mu = \mu'_1 = E[X]$.

Expected Value

- Recall that "average" is an arithmetic average where all available observations are weighted equally.
- The expected value, on the other hand, is the average of all possible values a random variable can take, weighted by the probability distribution.
- The question is, which value would we expect the random variable to take on, on average.

Definition 1.2

The expected value or mean of a random variable $g(X)$, denoted by $E[g(X)]$, is

$$E[g(X)] = \begin{cases} \int_{-\infty}^{\infty} g(x)f_X(x)dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} g(x)f_X(x) = \sum_{x \in \mathcal{X}} g(x)P(X = x) & \text{if } X \text{ is discrete} \end{cases}$$

If $E[g(X)] = \infty$, we say that $E[g(X)]$ does not exist.

- we are taking the average of $g(x)$ over all of its possible values ($x \in \mathcal{X}$), where these values are weighted by the respective value of the pdf, $f_X(x)$.

Example 1.1

Suppose X has an **exponential(λ) distribution**, that is, it has pdf given by

$$f_X(x) = \frac{1}{\lambda} e^{-x/\lambda}, \quad 0 \leq x < \infty \quad \lambda > 0.$$

Then,

$$E[X] = \int_0^{\infty} \frac{1}{\lambda} x e^{-x/\lambda} dx = -x e^{-x/\lambda} \Big|_0^{\infty} + \int_0^{\infty} e^{-x/\lambda} dx \quad (1)$$

$$= \int_0^{\infty} e^{-2/\lambda} dx = \lambda. \quad (2)$$

- To obtain this result, we use a method called integration by parts. This is based on

$$\int u dv = uv - \int v du.$$

- A very useful property of the expectation operator is that it is a linear operator.
- For example, consider some X such that $E[X] = \mu$.
- Then, for two constants a and b ,

$$E[a + Xb] = a + E[Xb] = a + bE[x] = a + b\mu.$$

- Notice that, clearly, the expectation of a constant is equal to itself.

Theorem 1.1

Let X be a random variable and let a , b and c be constants. Then for any functions $g_1(x)$ and $g_2(x)$ whose expectations exist,

- $E[ag_1(X) + bg_2(X) + c] = aE[g_1(X)] + bE[g_2(X)] + c.$
- If $g_1(x) \geq 0$ for all x , then $E[g_1(X)] \geq 0.$
- If $g_1(x) \geq g_2(x)$ for all x , then $E[g_1(X)] \geq E[g_2(X)].$
- If $a \leq g_1(x) \leq b$ for all x , then $a \leq E[g(X)] \leq b.$

Proof: Exercise!

Example 1.2

Let X have a uniform distribution, such that

$$f_X(x) = \begin{cases} 1 & \text{if } 0 \leq x \leq 1 \\ 0 & \text{if otherwise} \end{cases}$$

Define $g(X) = -\log X$. Then,

$$E[g(X)] = E[-\log X] = \int_0^1 -\log x dx = (-x \log x + x)|_0^1 = 1,$$

where we use integration by parts.

Variance

- variance measures the variation/dispersion/spread of the random variable around expectation.
- While the expectation is usually denoted by μ , σ^2 is generally used for variance.
- Variance is a second-order moment.

Definition 1.3

The variance of a random variable X is its **second central moment**,

$$\text{Var}(X) = E[(X - \mu)^2],$$

while $\sqrt{\text{Var}(X)}$ is known as the standard deviation of X .

- Importantly,

$$\text{Var}(X) = E[(X - \mu)^2] = E[X^2] - \mu^2.$$

- When it exists, the covariance of two random variables X and Y is defined as

$$\text{Cov}(X, Y) = E(\{X - E[X]\}\{Y - E[Y]\}).$$

Covariance and Correlation

- Let X and Y be two random variables. To keep notation concise, we will use the following notation.

$$E[X] = \mu_X, \quad E[Y] = \mu_Y, \quad \text{Var}(X) = \sigma_X^2 \quad \text{and} \quad \text{Var}(Y) = \sigma_Y^2.$$

Definition 2.1

The **covariance** of X and Y is the number defined by

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)].$$

Definition 2.2

The **correlation** of X and Y is the number defined by

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y},$$

which is also called the **correlation coefficient**.

Covariance and Correlation

- If **large**(**small**) values of X tend to be observed with **large**(**small**) values of Y , then will be positive.
- Why so? Within the above setting, when $X > \mu_X$ then $Y > \mu_Y$ is likely to be true whereas when $X < \mu_X$ then $Y < \mu_Y$ is likely to be true. Hence

$$E[(X - \mu_X)(Y - \mu_Y)] > 0.$$

- Similarly, if **large**(**small**) values of X tend to be observed with **small**(**large**) values of Y , then $\text{Cov}(X, Y)$ will be negative.

Covariance and Correlation

- Correlation normalises covariance by the standard deviations and is, therefore, a more informative measure.
- If $\text{Cov}(X, Y)=50$ while $\text{Cov}(W, Z)=0.9$, this does not necessarily mean that there is a much stronger relationship between X and Y . For example, if $\text{Var}(X)=\text{Var}(Y)=100$ while $\text{Var}(W)=\text{Var}(Z)=1$, then

$$\rho_{XY} = 0.5 \quad \rho_{WZ} = 0.9.$$

Theorem 2.1

For any random variables X and Y ,

$$\text{Cov}(X, Y) = E[XY] - \mu_X \mu_Y.$$

- Proof:**

$$\begin{aligned}\text{Cov}(X, Y) &= E[(X - \mu_X)(Y - \mu_Y)] \\ &= E[XY - \mu_X Y - \mu_Y X + \mu_X \mu_Y] \\ &= E[XY] - \mu_X E[Y] - \mu_Y E[X] + \mu_X \mu_Y \\ &= E[XY] - \mu_X \mu_Y.\end{aligned}$$

Theorem 2.2

If $X \perp\!\!\!\perp Y$, then $\text{Cov}(X, Y) = \rho_{XY} = 0$.

- **Proof:** Since $X \perp\!\!\!\perp Y$, by Theorem (2.1), Then

$$\text{Cov}(X, Y) = E[XY] - \mu_X \mu_Y = \mu_X \mu_Y - \mu_X \mu_Y = 0,$$

and consequently,

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = 0.$$

- It is crucial to note that although $X \perp\!\!\!\perp Y$ implies that $\text{Cov}(X, Y) = \rho_{XY} = 0$, the relationship does not necessarily hold in the reverse direction.

Theorem 2.3

If X and Y are any two random variables and a and b are any two constants, then

$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$$

If X and Y are independent random variables, then

$$\text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y).$$

- **Proof:** Exercise!

Covariance and Correlation

- Note that if two random variables, X and Y , are positively correlated, then

$$\text{Var}(X + Y) > \text{Var}(X) + \text{Var}(Y),$$

whereas if X and Y are negatively correlated, then

$$\text{Var}(X + Y) < \text{Var}(X) + \text{Var}(Y).$$

- For positively correlated random variables, large values in one tend to be accompanied by large values in the other. Therefore, the total variance is magnified.
- Similarly, for negatively correlated random variables, large values in one tend to be accompanied by small values in the other. Hence, the variance of the sum is dampened.

Variance of Sums of Random Variables

- Let a_i be some constant and X_i be some random variable, where $i = 1, \dots, n$.
- Then

$$\text{Var}\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n a_i^2 \text{Var}(X_i) + \sum_{i \neq j} \sum a_i a_j \text{Cov}(X_i, X_j).$$

Notice that $\sum_{i \neq j} \sum a_i a_j \text{Cov}(X_i, X_j)$ includes all possible covariances since it includes all the terms of the form $\text{Cov}(X_i, X_j)$ for all possible combinations of i and j such that $i \neq j$.

- An equivalent way of writing this is

$$2 \sum_{i < j} \sum a_i a_j \text{Cov}(X_i, X_j).$$

Variance of Sums of Random Variables

- Note that this second representation contains terms such as

$$\text{Cov}(X_1, X_2), \quad \text{Cov}(X_1, X_6), \quad \text{Cov}(X_{17}, X_{256}) \quad \text{etc.}$$

but NOT

$$\text{Cov}(X_2, X_1), \quad \text{Cov}(X_6, X_1), \quad \text{Cov}(X_{256}, X_{17}) \quad \text{etc.}$$

Hence, the double summation is multiplied by 2.

third and fourth moments

- third and fourth moments are concerned with how symmetric and fat-tailed the underlying distribution is.

Moment Generating Functions

- **moment generating function** can be used to obtain moments of a random variable.

Moments and Moment Generating Functions

Definition 3.1

Let X be a random variable with cdf F_X . The **moment generating function (mgf)** of X (or F_X), denoted by $M_X(t)$, is

$$M_X(t) = E[e^{tX}],$$

provided that the expectation exists for t in some neighbourhood of 0. That is, there is an $h > 0$ such that, for all t in $-h < t < h$, $E[e^{tX}]$ exists. If the expectation does not exist in a neighbourhood of 0, we say that the mgf does not exist.

- We can write the mgf of X as

$$M_X(t) = \int_{-\infty}^{\infty} e^{tx} f_X(x) dx \quad \text{if } X \text{ is continuous,}$$

$$M_X(t) = \sum_x e^{tx} P(X = x) \quad \text{if } X \text{ is discrete.}$$

Theorem 3.1

If X has mgf $M_X(t)$, then

$$E[X^n] = M_X^{(n)}(0),$$

where we define

$$M_X^{(n)}(0) = \frac{d^n}{dt^n} M_X(t)|_{t=0}.$$

That is, the n^{th} moment is equal to the n^{th} derivative of $M_X(t)$ evaluated at $t=0$.

- Now consider the pdf for $X \sim N(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right], \quad -\infty < x < \infty.$$

- The mgf is given by

$$M_X(t) = E[e^{Xt}] = \int_{-\infty}^{\infty} e^{tx} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2} dx.$$

- Then

$$M_X(t) = \exp(\mu t + \frac{\sigma^2 t^2}{2}).$$

- Clearly,

$$E[X] = \frac{d}{dt} M_X(t)|_{t=0} = (\mu + \sigma^2 t) \exp(\mu t + \frac{\sigma^2 t^2}{2})|_{t=0} = \mu,$$

$$\begin{aligned} E[X^2] &= \frac{d^2}{dt^2} M_X(t)|_{t=0} = \sigma^2 \exp(\mu t + \frac{\sigma^2 t^2}{2})|_{t=0} \\ &\quad + (\mu + \sigma^2 t)^2 \exp(\mu t + \frac{\sigma^2 t^2}{2})|_{t=0} \\ &= \sigma^2 + \mu^2, \end{aligned}$$

$$\text{Var}(X) = E[X^2] - \{E[X]\}^2 = \sigma^2 + \mu^2 - \mu^2 = \sigma^2.$$

Matrix Notation for Moments

- Let a_i, b_i be some numbers and X_i and Y_i be some random variables, where $i = 1, \dots, n$.
- Most of the time, we have to deal with terms such as

$$E\left[\sum_{i=1}^n a_i X_i\right] = \sum_{i=1}^n a_i E[X_i],$$

$$\text{Cov}\left(\sum_{i=1}^n a_i X_i, \sum_{j=1}^n b_j Y_j\right) = \sum_{i=1}^n \sum_{j=1}^n a_i b_j \text{Cov}(X_i, Y_j),$$

$$\text{Var}\left(\sum_{i=1}^n a_i X_i\right) = \sum_{i=1}^n \text{Var}(a_i X_i) + 2 \sum_{i < j} \sum a_i a_j \text{Cov}(X_i, X_j),$$

for different combinations of $i, j = 1, \dots, n$.

- The idea is to use matrix notation to represent such information more compactly and to manipulate it more easily.

Matrix Notation for Moments

- Our main object is some (random) vector,

$$X = (X_1, \dots, X_n)'.$$

Note that when one writes $X = (X_1, \dots, X_n)$, one means a row vector.

- Then,

$$E[X] = (E[X_1], \dots, E[X_n])'.$$

- Now, for $a = (a_1, \dots, a_n)'$,

$$E[a'X] = E\left[\sum_{i=1}^n a_i X_i\right] = \sum_{i=1}^n a_i E[X_i] = a' E[X].$$

Matrix Notation for Moments

- Now, define an $(n \times k)$ matrix

$$U = \begin{bmatrix} U_{11} & \dots & U_{1k} \\ \vdots & & \vdots \\ U_{n1} & \dots & U_{nk} \end{bmatrix}, \text{ where } E[U] = \begin{bmatrix} E[U_{11}] & \dots & E[U_{1k}] \\ \vdots & & \vdots \\ E[U_{n1}] & \dots & E[U_{nk}] \end{bmatrix}$$

where U_{ij} is the entry for row i and column j , $i = 1, \dots, n$ and $j = 1, \dots, k$.

- Same as before,

$$E[a'U] = a'E[U], \quad E[Ub] = E[U]b, \quad \text{and} \quad E[a'Ub] = a'E[U]b.$$

- For example,

$$E[a'Ub] = E\left[\sum_{i=1}^n \sum_{j=1}^k a_i b_j U_{ij}\right] = \sum_{i=1}^n \sum_{j=1}^k a_i b_j E[U_{ij}] = a'E[U]b$$

Matrix Notation for Moments

- Now, let X and Y be $(r \times 1)$ and $(c \times 1)$ random vectors, respectively. Define
- In other words,

$$\begin{aligned} \text{Cov}(X, Y) &= \begin{bmatrix} \text{Cov}(X_1, Y_1) & \dots & \text{Cov}(X_1, Y_c) \\ \vdots & & \vdots \\ \text{Cov}(X_r, Y_1) & \dots & \text{Cov}(X_r, Y_c) \end{bmatrix} \\ &= E \begin{bmatrix} \{X_1 - E[X_1]\}\{Y_1 - E[Y_1]\} & \dots & \{X_1 - E[X_1]\}\{Y_c - E[Y_c]\} \\ \vdots & & \vdots \\ \{X_r - E[X_r]\}\{Y_1 - E[Y_1]\} & \dots & \{X_r - E[X_r]\}\{Y_c - E[Y_c]\} \end{bmatrix} \end{aligned}$$

Matrix Notation for Moments

$$\begin{aligned} &= E \left[\begin{pmatrix} X_1 - E[X_1] \\ \vdots \\ X_r - E[X_r] \end{pmatrix} (Y_1 - E[Y_1], \dots, Y_c - E[Y_c]) \right], \\ &= E(\{X - E[X]\}\{Y - E[Y]\}'). \end{aligned}$$

Matrix Notation for Moments

- Usually, for a $(c * 1)$ vector X , one would write $Cov(X)$ for $Cov(X, X)$,
- This is given by

$$= Cov(X) \begin{bmatrix} Var(X_1) & \dots & Cov(X_1, X_c) \\ \vdots & & \vdots \\ Cov(X_1, X_c) & \dots & Var(X_c) \end{bmatrix},$$

which is a $(c * c)$ symmetric matrix.

Matrix Notation for Moments

- We can also consider block structures. Let

$$X = \begin{pmatrix} Y \\ Z \end{pmatrix},$$

where Y is $(p * 1)$ vector and Z is a $(q * 1)$ vector.

- Then,

$$\begin{aligned} \text{Cov}(X) &= E\left(\left\{\begin{pmatrix} Y \\ Z \end{pmatrix} - E\left[\begin{pmatrix} Y \\ Z \end{pmatrix}\right]\right\}\left\{\begin{pmatrix} Y \\ Z \end{pmatrix} - E\left[\begin{pmatrix} Y \\ Z \end{pmatrix}\right]\right\}'\right) \\ &= E\begin{pmatrix} \{Y - E[Y]\}\{Y - E[Y]\}' & \{Y - E[Y]\}\{Z - E[Z]\}' \\ \{Z - E[Z]\}\{Y - E[Y]\}' & \{Z - E[Z]\}\{Z - E[Z]\}' \end{pmatrix} \\ &= \begin{pmatrix} \text{Cov}(Y) & \text{Cov}(Y, Z) \\ \text{Cov}(Z, Y) & \text{Cov}(Z) \end{pmatrix}, \end{aligned}$$

where $\text{Cov}(Y)$ is $(p * p)$, $\text{Cov}(Y, Z)$ is $(p * q)$, $\text{Cov}(Z, Y)$ is $(q * p)$ and $\text{Cov}(Z)$ is $(q * q)$.

Matrix Notation for Moments

- For such block structures, the following result might one day come in very handy.
- Let A_{11} be an $(m_1 * m_1)$, A_{12} be an $(m_1 * n_2)$, A_{21} be an $(m_2 * n_1)$ and A_{22} be an $(m_2 * n_2)$ matrix.
- Then,

$$\begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}' = \begin{pmatrix} A^{11} & A^{12} \\ A^{21} & A^{22} \end{pmatrix},$$

where

$$A^{11} = A_{11}^{-1} + A_{11}^{-1} A_{12} D^{-1} A_{21} A_{11}^{-1},$$

$$A^{12} = -A_{11}^{-1} A_{12} D^{-1},$$

$$A^{21} = -D^{-1} A_{21} A_{11}^{-1},$$

$$A^{22} = D^{-1},$$

$$D = A_{22} - A_{21} A_{11}^{-1} A_{12},$$

as long as all the inverses exist.

Matrix Notation for Moments

- The following alternative would also work:

$$A^{11} = E^{-1},$$

$$A^{12} = -E^{-1}A_{12}A_{22}^{-1},$$

$$A^{21} = -A_{22}^{-1}A_{21}E^{-1},$$

$$A^{22} = A_{22}^{-1} + A_{22}^{-1}A_{21}E^{-1}A_{12}A_{22}^{-1},$$

$$E = A_{11} - A_{12}A_{22}^{-1}A_{21}.$$

Matrix Notation for Moments

- Let a and b be $(r \times 1)$ and $(c \times 1)$ non-stochastic vectors. We might encounter terms such as $\text{Cov}(a'X, b'Y)$ or $\text{Var}(a'X)$.
- Let $E[X_i] = \mu_{X_i}$, $E[Y_i] = \mu_{Y_i}$ and $\text{Cov}(X_i, Y_j) = \sum_{X_i, Y_j}$. Then

$$\begin{aligned}\text{Cov}(a'X, b'Y) &= \text{Cov}\left(\sum_{i=1}^r a_i X_i, \sum_{j=1}^c b_j Y_j\right) \\&= E\left\{\left[\sum_{i=1}^r a_i (X_i - \mu_{X_i})\right]\left[\sum_{j=1}^c b_j (Y_j - \mu_{Y_j})\right]\right\} \\&= \sum_{i=1}^r \sum_{j=1}^c a_i b_j E[(X_i - \mu_{X_i})(Y_j - \mu_{Y_j})] \\&= \sum_{i=1}^r \sum_{j=1}^c a_i b_j \sum_{X_i, Y_j} = a' \sum_{XY} b = a' \text{Cov}(X, Y) b.\end{aligned}$$

Matrix Notation for Moments

- Now, let $\Sigma_{ij} = \text{Cov}(X_i, X_j)$ and $\Sigma_{XX} = \text{Var}(X)$. Then,

$$\begin{aligned}\text{Var}(a'X) &= E[(\sum_{i=1}^r a_i X_i - E[\sum_{i=1}^r a_i X_i])^2] \\&= E\{[\sum_{i=1}^r a_i (X_i - \mu_i)][\sum_{i=1}^r a_i (X_i - \mu_i)]\} \\&= \sum_{i=1}^r \sum_{j=1}^r a_i a_j E[(X_i - \mu_i)(X_j - \mu_j)] \\&= \sum_{i=1}^r \sum_{j=1}^r a_i a_j \Sigma_{ij} = a' \text{Var}(X) a.\end{aligned}$$

Matrix Notation for Moments

- Now, Consider

$$\begin{aligned} \text{Var}(X + Y) &= E\{[(X - \mu_X) + (Y - \mu_Y)][(X - \mu_X) + (Y - \mu_Y)]'\} \\ &= E[(X - \mu_X)(X - \mu_X)'] + E[(X - \mu_X) + (Y - \mu_Y)]' \\ &\quad + E[(Y - \mu_Y)(X - \mu_X)'] + E[(Y - \mu_Y) + (Y - \mu_Y)]' \\ &= \Sigma_{XX} + \Sigma_{XY} + \Sigma_{YX} + \Sigma_{YY}. \end{aligned}$$

- Using this, we get

$$\begin{aligned} \text{Var}[a'(X + Y)] &= a'(\Sigma_{XX} + \Sigma_{XY} + \Sigma_{YX} + \Sigma_{YY})a \\ &= a'\Sigma_{XX}a + 2a'\Sigma_{XY}a + a'\Sigma_{YY}a, \end{aligned}$$

where we use the fact that

$$a'\Sigma_{XY}a = a'\Sigma_{YX}a$$

- These results easily extend to cases where a and b are replaced by matrices.

$$E[RX] = RE[X]$$

$$\begin{aligned} \text{Var}(RX) &= E[R(X - \mu_X)(X - \mu_X)'R'] \\ &= RE[(X - \mu_X)(X - \mu_X)']R' \\ &= R\Sigma_{XX}R'. \end{aligned}$$

Discrete Uniform Distribution

- A random variable X has a **discrete uniform**(1, N) distribution if

$$P(X = x|N) = \frac{1}{N}, \quad x = 1, 2, \dots, N,$$

where N is a specified integer. This distribution puts equal mass on each of the outcomes $1, 2, \dots, N$.

Binomial Distribution

- This is based on a Bernoulli trial (after James Bernoulli which is an experiment with two, and only, two, possible outcomes.
- A random variable X has Bernoulli(p) distribution if

$$X = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1 - p \end{cases} \quad 0 \leq p \leq 1.$$

- $X = 1$ is often termed as "success" and p is, accordingly, the probability of success. Similarly, $X = 0$ is termed a "failure".
- Now,

$$E[X] = 1 * p + 0 * (1 - p) = p,$$

$$\text{and } \text{Var}(X) = (1 - p)^2 p + (0 - p)^2 (1 - p) = p(1 - p).$$

- $E[X] = np$ (**Proof:** Exercise!)
- $\text{Var}(X) = np(1 - p)$ (**Proof:** Exercise!)
- Examples:
 - 1 Tossing a coin (p = probability of a head, $X = 1$ if heads)
 - 2 Roulette ($X = 1$ if red occurs, p = probability of red)
 - 3 Election polls ($X = 1$ if candidate A gets vote)
 - 4 Incidence of disease (p = probability that a random person gets infected)

Binomial Distribution

- We can extend the scope to a collection of many independent trials.
- Define

$$A_i = \{X = 1 \text{ on the } i^{\text{th}} \text{ trial}\}, \quad i = 1, 2, \dots, n.$$

- Assuming that A_1, \dots, A_n are independent events, we can derive the **distribution of the total number of successes in n trials**. Define $Y =$ "total number of successes in n trials".
- The event $\{Y = y\}$ means that out of n trials, y resulted as success. Therefore, $n - y$ trials have been unsuccessful.
- In other words, exactly y of A_1, \dots, A_n must have occurred.
- There are many possible orderings of the events that would lead to this outcome. Any particular such ordering has probability

$$p^y(1 - p)^{n-y}.$$

- Since there are $\binom{n}{y}$ such sequences, we have

$$P(Y = y | n, p) = \binom{n}{y} p^y (1 - p)^{n-y}, \quad y = 0, 1, 2, \dots, n,$$

and Y is called a *binomial*(n, p) random variable.

Example 5.1

If you flip a fair coin 10 times, what is the probability of getting all tails?

- Let's first calculate the probability of getting tail on fair coin when you flip it one time.

$$P(1) = \frac{1}{2} = 50\% \text{ (Because coin has two sides, H \& T)}$$

- Since all the trails are independent, probability of getting head on n th turn is also $1/2$.
- Then,

$$\begin{aligned} P(10) &= \frac{1}{2} * \frac{1}{2} * \dots * \frac{1}{2} \quad (10 \text{ times}) \\ &= \left(\frac{1}{2}\right)^{10}. \end{aligned}$$

Poisson Distribution

- In modelling a phenomenon in which we are waiting for an occurrence (such as waiting for a bus), **the number of occurrence in a given time interval** can be modelled by the Poisson distribution.
- The basic assumption is as follows: for small time intervals, the probability of an arrival is proportional to the length of waiting time.
- If we are waiting for the bus, the probability that a bus will arrive within the next hour is higher than the probability that it will arrive within 5 minutes.
- Other possible applications are distribution of bomb hits in an area or distribution of fish in a lake.
- The only parameter is λ , also sometimes called the "intensity parameter."

Poisson Distribution

- $E[X] = \lambda$
- $Var(X) = \lambda$
- **Proof:** Exercise!

Example 5.2

As an example of a waiting-for-occurrence application, consider a telephone operator who, on average, handles five calls every 3 minutes. What is the probability that there will be no calls in the next minute? At least two calls? If we let X = number of calls in a minute, then X has a Poisson distribution with $E[X] = \lambda = 5/3$. So,

$$P(\text{no calls in the next minute}) = P(X = 0)$$

$$= \frac{e^{-5/3}(5/3)^0}{0!} = e^{-5/3} = 0.189$$

$$\text{and} \quad P(\text{at least two calls in the next minute}) = P(X \geq 2)$$

$$= 1 - P(X = 0) - P(X = 1)$$

$$= 1 - 0.189 - \frac{e^{-5/3}(5/3)^1}{1!}$$

$$= 0.496.$$

Uniform Distribution

- The continuous uniform distribution is defined by spreading mass uniformly over an interval $[a, b]$. Its pdf is given by

$$f(x|a, b) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{if otherwise} \end{cases}.$$

- One can easily show that

$$\int_a^b f(x) dx = 1,$$

$$E[X] = \frac{b+a}{2},$$

$$\text{Var}(X) = \frac{(b-a)^2}{12}.$$

- In many cases, when people say Uniform distribution, they implicitly mean $(a, b) = (0, 1)$.

Exponential Distribution

- Now consider, $\alpha = 1$:

$$f(x|a, \beta) = f(x|1, \beta) = \frac{1}{\beta} e^{-x/\beta}, \quad 0 < x < \infty.$$

- Again, using our previous results, for X exponential (β) we have

$$E[X] = \beta \quad \text{and} \quad \text{Var}(X) = \beta^2$$

- A peculiar feature of this distribution is that **it has no memory**.

Exponential Distribution

- If $X \sim \text{exponential}(\beta)$, then, for $s > t \geq 0$,

$$\begin{aligned}P(X > s | X > t) &= \frac{P(X > s, X > t)}{P(X > t)} = \frac{P(X > s)}{P(X > t)} \\&= \frac{\int_s^\infty \frac{1}{\beta} e^{-x/\beta} dx}{\int_t^\infty \frac{1}{\beta} e^{-x/\beta} dx} = \frac{e^{-s/\beta}}{e^{-t/\beta}} \\&= e^{-(s-t)/\beta} = P(X > s - t).\end{aligned}$$

- This is because,

$$\int_{s-t}^\infty \frac{1}{\beta} e^{-x/\beta} dx = -e^{-x/\beta} \Big|_{x=s-t}^\infty = e^{-(s-t)/\beta}.$$

- What does this mean? When calculating $P(X > s | X > t)$, what matters is not whether X has passed a threshold or not. What matters is the distance between the threshold and the value to be reached.
- If Mr X has been fired more than 10 times, what is the probability that he will be fired more than 12 times? It is not different from the probability that a person, who has been fired once, will be fired more than two times. History does not matter.

Normal Distribution

- We now consider the **normal distribution** or the **Gaussian distribution**.
- Why is this distribution so popular?
 - 1 Analytical tractability
 - 2 Bell shaped or symmetric
 - 3 It is central to Central Limit Theorem; this type of results guarantee that, under (mild) conditions, the normal distribution can be used to approximate a large variety of distribution in large samples.
- The distribution has two parameters: mean and variance, denoted by μ and σ^2 , respectively.
- The pdf is given by,

$$f(x|\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-1/2 \frac{(x - \mu)^2}{\sigma^2}\right].$$

Normal Distribution

- This distribution is usually denoted as $N(\mu, \sigma^2)$.
- A very useful result is that for $X \sim N(\mu, \sigma^2)$,

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1).$$

- $N(0, 1)$ is known as the standard normal distribution.
- To see this, consider the following:

$$\begin{aligned} P(Z \leq z) &= P\left(\frac{(X - \mu)/\sigma \leq z}{1}\right) \\ &= P(X \leq z\sigma + \mu) \\ &= \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{z\sigma + \mu} e^{-(x-\mu)^2/2\sigma^2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-t^2/2} dt, \end{aligned}$$

where we substitute $t = (x - \mu)/\sigma$. Notice that this implies that $dt/dx = 1/\sigma$. This shows that $P(Z \leq z)$ is the standard normal cdf.

Normal Distribution

- Then, we can do all calculations for the standard normal variable and then convert these results for whatever normal random variable we have in mind.
- Consider, for $Z \sim N(0, 1)$, the following:

$$E[Z] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} ze^{-z^2/2} dz = -\frac{1}{\sqrt{2\pi}} e^{-z^2/2} \Big|_{-\infty}^{\infty} = 0.$$

- Then, to find $E[X]$ for $X \sim N(\mu, \sigma^2)$, we can use $X = \mu + Z\sigma$:

$$E[X] = E[\mu + Z\sigma] = \mu + \sigma E[Z] = \mu + \sigma * 0 = \mu.$$

- Similarly,

$$\text{Var}(X) = \text{Var}(\mu + Z\sigma) = \sigma^2 \text{Var}(Z) = \sigma^2.$$

- What about

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-z^2/2} dz \stackrel{?}{=} 1.$$

Lognormal Distribution

- Let X be a random variable such that

$$\log X \sim N(\mu, \sigma^2).$$

Then, X is said to have a lognormal distribution.

- By using a transformation argument (Theorem (1.2)), the pdf of X is given by,

$$f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \frac{1}{x} \exp\left[-\frac{(\log x - \mu)^2}{2\sigma^2}\right],$$

where $0 < x < \infty$, $-\infty < \mu < \infty$, and $\sigma > 0$.

- How? Take $W = \log X$. We start from distribution of W and want to find the distribution of $X = \exp W$. Then, $g(W) = \exp(W)$ and $g^{-1}(X) = \log(X)$. The rest follows by using Theorem (1.2).

Laplace distribution

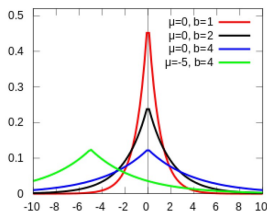
Laplace prior for β :

$$\beta_j \sim \text{Lap}(0, \frac{2\sigma^2}{\lambda}) \Rightarrow p(\beta_j) = \frac{\lambda}{4\sigma^2} \exp(-\frac{\lambda}{2\sigma^2} |\beta_j|)$$

where $\beta_j, j = 1, \dots, p$ are i.i.d

- If $Z \sim \text{Lap}(\mu, b)$, then $E[Z] = \mu$, $\text{Var}(Z) = 2b^2$,

$$p(Z) = \frac{1}{2b} \exp(-\frac{|x - \mu|}{b})$$



- **Likelihood:** Gaussian likelihood

Casella, G., & Berger, R. (2002). *Statistical inference*. Cengage Learning.
<https://books.google.fr/books?id=FAUVEAAAQBAJ>