

Introduction to Statistical Modelling

Semester 2, 2022

Dr Gentry White

Tarang Janawalkar

This work is licensed under a Creative Commons
“Attribution-NonCommercial-ShareAlike 4.0 International” license.



Contents

Contents	1
1 Introduction	2
1.1 Elements of Statistical Modelling	2
1.1.1 Data	2
1.1.2 Collecting information	2
1.1.3 Randomness	2
1.1.4 Probability	2
1.2 Experimental Units and Populations	2
1.2.1 Sample vs. Population	2
1.3 Types of Data	3
1.3.1 Univariate, Bivariate, and Multivariate	3
1.3.2 Experimental vs. Observational Data	3
1.3.3 Quantitative Data	3
1.3.4 Qualitative Data	3
1.4 Summarising and Describing Data	3
1.4.1 Tables	3
1.5 Bar Charts	4
1.6 Line Charts	4
1.7 Histograms	4
1.8 Plots, Graphs, and Charts	4
1.9 Interpreting Graphical Descriptions	4
1.9.1 Centrality	5
1.9.2 Skew	5
1.9.3 Trends	5
2 Numerical Summaries of Data	5
2.1 Measures of Centrality	5
2.1.1 Mean	5
2.1.2 Median	6
2.1.3 Mode	6
2.1.4 Population Mean	6
2.2 Measures of Dispersion	6
2.2.1 Range	6
2.2.2 Variance	6
2.2.3 Standard Deviation	7
2.3 Skew	7
2.4 Measures of Rank	8
2.4.1 Z-Score	8
2.4.2 Quantiles	8
2.4.3 Inter-Quartile Range	8
2.5 Boxplots	8
2.5.1 Five Number Summary	8
2.5.2 Outliers	9

3	Bivariate Data	9
3.1	Bivariate Categorical Data	9
3.1.1	Contingency Tables	9
3.1.2	Bar Plots	9
3.2	Bivariate Quantitative Data	9
3.3	Scatter Plots	9
3.3.1	Covariance and Correlation Coefficients	9
3.4	Regression and Least Squares	10
4	Probability	10
4.1	Experiments, Events, Sample Space	10
4.2	Probability of Events	11
4.2.1	Probability of an Event	11
4.2.2	Probability of an Event in a Sample Space	11
4.2.3	Probability of the Complement	11
4.2.4	Probability of Subsets	12
4.2.5	Addition Law	12
4.3	Conditional Probability	12
4.3.1	Independence	12
4.4	Bayes' Rules	12
4.4.1	Law of Total Probability	12
5	Probability Distributions	13
5.1	Random Variables	13
5.2	Discrete Random Variables	13
5.2.1	Probability Mass Function	13
5.2.2	Cumulative Mass Function	13
5.2.3	Expectation	14
5.2.4	Median and Mode	14
5.2.5	Variance	14
5.3	Continuous Random Variables	14
5.3.1	Probability Density Function	14
5.3.2	Cumulative Density Function	15
5.3.3	Median and Mode	15
5.3.4	Expectation	15
5.3.5	Variance	15
5.4	Probability Distributions	16
5.4.1	Bernoulli Distribution	16
5.4.2	Binomial Distribution	16
5.4.3	Poisson Distribution	17
5.4.4	Uniform Distribution	18
5.4.5	Exponential Distribution	18
5.4.6	Memoryless Property	19
5.4.7	Normal Distribution	19
5.5	Standard Normal Distribution	19

6	Sampling	20
6.1	Observational and Experimental Studies	20
6.1.1	Observational Studies	20
6.1.2	Experimental Studies	20
6.2	Sampling	20
6.2.1	Simple Random Sampling	20
6.3	Stratified Random Sampling	20
6.4	Cluster Sampling	20
6.5	Non-random Sampling Methods	21
7	Sampling Distributions	21
7.1	Central Limit Theorem	21
7.2	Standard Error	21
7.3	Sample Proportion	21
7.4	Assessing Normality	22

1 Introduction

Statistics is a field of mathematics that deals with data. It includes the study of summarising data, constructing probabilistic models, estimating parameters, and making statistical inferences. Statistical modelling includes asking questions, obtaining data and determining a mathematical model.

1.1 Elements of Statistical Modelling

1.1.1 Data

Data is a collection of numbers that describes some characteristic that can be ranked, counted, or measured.

1.1.2 Collecting information

Statistical modelling relies upon reliably sourced data. When collecting data, we must consider

- what questions are we trying to answer,
- what information is needed to answer these questions,
- what is the best source for that information

1.1.3 Randomness

We must be aware that everything is different and that randomness introduces uncertainty in data. Random events are events whose exact outcome cannot be predicted. We can assume that all variation in the world is observed due to randomness.

1.1.4 Probability

Probability is a mathematical construct for dealing with randomness and uncertainty.

1.2 Experimental Units and Populations

Definition 1.1 (Experimental unit). An **experimental unit** is an individual that generates information for the data collection process. Careful consideration of what constitutes an experimental unit must be made to ensure that it aligns with the questions of interest.

1.2.1 Sample vs. Population

Definition 1.2 (Population). We might have questions about a very large collection of things called a **population**.

A dataset collected from a population is called a census.

As it is not feasible to collect data from an entire population, we must use a sample of the population.

Definition 1.3 (Sample). A **sample** is a subset of a population that is representative of the population, in some cases a random sample is sufficient.

Definition 1.4 (Random sample). A **random sample** is one where the sample members are selected from the population by chance.

1.3 Types of Data

1.3.1 Univariate, Bivariate, and Multivariate

Data can be described in terms of dimension, that is, how many measurements were collected from each experimental unit. By collecting multiple measurements from each experimental unit, we can ask questions about the relationship between the measurements.

- When a single measurement is collected, the resulting dataset is **univariate**.
- If two measurements are collected, the dataset is **bivariate**.
- If more than two measurements are collected, the dataset is **multivariate**.

1.3.2 Experimental vs. Observational Data

Data sets that have been collected without any specific analyses or modelling in mind are called **observational data**. By contrast, when a collection procedure is specifically designed to obtain data with a specific intent, i.e., a laboratory test, the data is called **experimental data**. Observational data may contain biases that limit its usefulness and bias any modelling or analysis results.

1.3.3 Quantitative Data

Quantitative data is data that is expressed numerically. This data can be classified as *discrete*, *continuous*, or *ordinal*.

- Count data is classified as discrete, i.e., integer values or finite sets of real values.
- Continuous data is a measurement on a continuum or a measure that can be subdivided infinitely, i.e., time and lengths.
- Ordinal data is data where the order or ranking of values (discrete or continuous) is important.

When data is not ordinal, it is called **nominal** data.

1.3.4 Qualitative Data

Qualitative (categorical) data is data where the variable of interest is membership to a group or category.

1.4 Summarising and Describing Data

1.4.1 Tables

Tables are the most immediate way of summarising a data set. We might organise data in a table with one row for each subject and a column for each measurement.

1.5 Bar Charts

Graphical depictions of the data can also be useful but are limited in the number of variables displayed in one picture.

Bar charts are most useful for categorical data where categories are listed on the x -axis of the plot, and bars for each category are drawn with their heights corresponding to the *counts* for that category.

When the categories are **ordered** from left to right in descending order counts, the plot is called a **Pareto plot**.

1.6 Line Charts

Line charts illustrate a *trend* of change based on **two** quantitative variables. Typically line charts display trends over time (or other ordinal variables).

Often trends over time need to be aggregated by plotting the average or median per year to avoid a “busy” plot which can sometimes be difficult to read.

While the resulting chart can explain overall trends, they can obscure how much variability or “noise” is in the data and may be misleading if the overall trend is obscured by variability.

1.7 Histograms

Histograms are a special kind of bar chart that give a visual description of data by “binning” or grouping data into data ranges, then plotting bars with heights equal to the count of the bins’ contents *or* the relative proportion of the bins’ contents.

Histograms give us a picture of the shape of the data and help identify patterns in the distribution of values.

The binning process is performed by the computer, however in most cases we override the automatic settings and select either the number of bins, or the width of each bin.

1.8 Plots, Graphs, and Charts

- A **chart** is a visual display of data, i.e., a table, a graph, or a diagram
- A **graph** is a diagram showing the relationship between variables, each measured along orthogonal axes.
- A **plot** is used as a synonym for graph but is less precise in its definition; it also sometimes refers specifically to a graph *produced by a computer*.

1.9 Interpreting Graphical Descriptions

Graphical descriptions of data should ensure that all information about the data is expressed.

- The x and y axes should be clear in what they are measuring, including any units.
- Consider how the graph or chart was made. What choices were made and how might different options change how the graph is perceived.

- Does the graph contain any outliers that merit investigation to determine if they are accurate measurements, or if they result from either measurement or recording error.
- For Pareto charts and histograms; the y -axis should measure proportion or density rather than frequency to make comparisons easier.

1.9.1 Centrality

Histograms are a graphical representation of the distribution or density of observations. Centrality is the degree to which an observation is central to the distribution. Additionally, the data can be multi-modal if there are multiple “peaks” or “centres” in the distribution.

Altering the number of bins or bin width may reveal the centrality of the observations.

1.9.2 Skew

Another characteristic of histograms is the degree to which the distribution is skewed. Skew is the deviation from symmetry about the centre of the data. Skew is either “right” skew where the tail of the density or histogram is heavier to the right, or “left” skew if otherwise.

This can be observed by looking at how much the left/right tails are stretched in comparison to one another, i.e., the tail to the right of a right skewed chart stretches further on the x -axis than on the left.

1.9.3 Trends

Trends refer to changes in a line chart and are often described as a constant (first-derivative) pattern of increasing or decreasing values.

2 Numerical Summaries of Data

Although graphical summaries are useful for developing a general understanding data, they are limited to subjective interpretations. To form a precise understanding of the data, we need to use numerical summaries. Here we must make a distinction between sample and population summaries as measurements may vary between samples, whereas population summaries are generally constant.

2.1 Measures of Centrality

2.1.1 Mean

Given a set of n observations x_1, x_2, \dots, x_n , the **arithmetic mean** or **average** is defined as

$$\frac{1}{n} \sum_{i=1}^n x_i \equiv \frac{1}{n} (x_1 + x_2 + \dots + x_n)$$

If the data is taken from a sample, the sample mean is denoted \bar{x} .

2.1.2 Median

A drawback to the mean is that it can be misleading when the data is skewed. The **median** is the middle value of a set of n observations when arranged from largest to smallest.

If n is odd:

$$\text{median} = x^{(\frac{n+1}{2})}$$

or the $(n+1)/2$ th value of the sorted list. If n is even, the median is the :

$$\text{median} = \frac{x^{(\frac{n}{2})} + x^{(\frac{n}{2}+1)}}{2}$$

2.1.3 Mode

Given discrete data, the mode is defined as the most common value in a set of observations.

2.1.4 Population Mean

The mean of a finite population is computed in the same way as the mean of a sample, but the population mean is denoted by μ .

2.2 Measures of Dispersion

Dispersion refers to how much variation there is in a set of observations.

2.2.1 Range

Given a set of observations that are ordered such that

$$x^{(1)} \leq x^{(2)} \leq \dots \leq x^{(n)}$$

the range is defined as

$$x^{(n)} - x^{(1)}.$$

2.2.2 Variance

The variance is the average of the squared deviations from the mean.

- Given the observations x_1, x_2, \dots, x_N , from a population of size N with mean μ , the **population variance** is defined as

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2.$$

- Given the observations x_1, x_2, \dots, x_n , from a sample of size n with mean \bar{x} , the **sample variance** is defined as

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

The **population variance** is given by

2.2.3 Standard Deviation

The standard deviation is the square root of the variance. This is conceptually easier to understand as it has the same units as the data.

- The **population standard deviation** is defined as

$$\sigma = \sqrt{\sigma^2}.$$

- The **sample standard deviation** is defined as

$$s = \sqrt{s^2}.$$

Theorem 2.2.1 (Chebyshev's Theorem). *Given a set of n observations, at least*

$$1 - \frac{1}{k^2}$$

of them are within k standard deviations of the mean, where $k \geq 1$. Formally,

$$\frac{\#\{x | \bar{x} - ks < x < \bar{x} + ks\}}{n} \geq 1 - \frac{1}{k^2}$$

Theorem 2.2.2 (Empirical Rule). *If a histogram of the data is approximately unimodal and symmetric, then,*

- 68% of the data falls within **one** standard deviation of the mean
- 95% of the data falls within **two** standard deviations of the mean
- 99% of the data falls within **three** standard deviations of the mean

Often the standard deviation cannot be computed directly, but can be approximated using the Empirical rule. Here we assume that

$$\text{range} \approx 4s$$

so that

$$s = \frac{\text{range}}{4}.$$

2.3 Skew

The **skew** describes the asymmetry of the distribution. For a finite population of size N , the **population skew** is defined as

$$\frac{1}{N} \sum_{i=1}^N \left(\frac{x_i - \mu}{\sigma} \right)^3$$

For a sample of size n , the **sample skew** is defined as

$$\frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\left(\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \right)^3}$$

- When the skew is **positive**, the data is **right-skewed** and the “tail” of the distribution is **longer on the right**
- When the skew is **negative**, the data is **left-skewed** and the “tail” of the distribution is **longer on the left**

2.4 Measures of Rank

It is often useful to know the rank or *relative standing* of a value in a set of observations. This is natural for ordinal data whose ordering has implicit meaning, but it can also be useful for nominal data as a means of measuring dispersion.

2.4.1 Z-Score

The Z-score is a unitless quantity and can be used to make comparisons of relative rank between members of a population.

$$Z = \frac{x - \mu}{\sigma} \quad \text{or} \quad \frac{x - \bar{x}}{s}$$

2.4.2 Quantiles

In addition to Z-scores, quantiles can be used to make comparisons of relative ranking between populations, as well as construct intervals bounding a given proportion of the observations. For a set of n observations, x_q is the q -th quantile, if $q\%$ of the observations are less than x_q .

2.4.3 Inter-Quartile Range

The inter-quartile range (IQR) is the difference between the 75th and 25th quantiles, or the range covered by the middle 50% of data. It is a robust measure of the dispersion of the data, as it is not affected by extreme values unlike the range or variance.

2.5 Boxplots

2.5.1 Five Number Summary

The five number summary is set of measurements that indicates the

- minimum value
- 25% quartile
- median
- 75% quartile
- maximum value

A boxplot is a graphical display of the five number summary. It is a plot of the values of the data mapped to the y -axis.

Using the `ggplot2` package, the function `geom_boxplot()` draws a box encompassing the IQR with a horizontal line indicating the median. Vertical lines extend 1.5 times the IQR above and below the box. The points not within the ends of the vertical lines are also plotted to indicate outliers.

2.5.2 Outliers

Outliers are extreme observations that fall outside some interval defined either by quantiles (above 95% or below 5% quantiles) or in terms of the Empirical rule (outside two standard deviations from the mean). They should be investigated to determine if they are errors or naturally occurring extreme values.

3 Bivariate Data

Data in two dimensions is often used to describe relationships between two variables.

3.1 Bivariate Categorical Data

Bivariate categorical data is a dataset with two qualitative or categorical variables that have a relationship we want to summarise. This can be done using contingency tables or side-by-side (or stacked) bar charts.

3.1.1 Contingency Tables

Contingency tables or crosstabs are tabular representations of the frequency of occurrence of pairs of values. The categories for each variable are assigned to an axis of the table so that each cell represents the frequency of occurrence of a pair of categories, one from each variable.

3.1.2 Bar Plots

Often the data is presented more effectively as a stacked bar chart or side-by-side bar chart. Here the counts for each pair of categories are plotted on the same axis, and stacked on top of one another to display relative proportion, or side-by-side if too busy.

3.2 Bivariate Quantitative Data

Bivariate quantitative data is a dataset with one qualitative variable and one quantitative variable. This can be represented as a table, or through various charts, by comparing charts side-by-side for each category.

3.3 Scatter Plots

When both variables are quantitative, the data can be represented as a scatter plot with each variable assigned to an axis and the points plotted on the axes.

3.3.1 Covariance and Correlation Coefficients

For such data, the covariance is the measure of the linear correlation between the variables. For variables x and y ,

$$s_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}.$$

Note that when $x = y$, the formula simplifies to the sample variance of x . The covariance has the following characteristics:

- $s_{xy} > 0$: As x increases, y also increases.
- $s_{xy} < 0$: As x increases, y decreases.
- $s_{xy} \approx 0$: No relationship between x and y .

Although the covariance is a useful tool to measure relationships, it is only generalisable in terms of its sign. Thus, if we want to compare across data sets, we need to use the correlation coefficient.

$$r_{xy} = \frac{s_{xy}}{s_x s_y}$$

The **correlation coefficient** is a measure of the strength of the relationship between the variables. It is a scale-free and unitless measure bounded between -1 and 1 and has the same characteristics as the covariance.

Note that a correlation coefficient of 0 indicates **no linear relationship** between the variables, and not necessarily indicative of **no relationship**.

3.4 Regression and Least Squares

In addition to the numerical summaries above, a regression or least squares line of best fit provides both a graphical and numerical summary of the relationship between the variables. A linear relationship between two variables x and y is defined as

$$y = a + bx.$$

The least squares best fit determines the coefficients a and b that minimise the sum of the squares of the residuals (errors) between y and the line $\hat{y} = a + bx$. Mathematically,

$$\min_{a, b} \sum_{i=1}^n (y_i - \hat{y}(x_i))^2.$$

The coefficients can be summarised by the formula

$$b = r \frac{s_y}{s_x} = \frac{s_{xy}}{s_x^2}$$

$$a = \bar{y} - b\bar{x}.$$

4 Probability

4.1 Experiments, Events, Sample Space

Definition 4.1 (Experiment). An experiment is a situation that produces some observable phenomena where the outcome is impossible to predict with certainty.

Definition 4.2 (Simple event). A simple event is the outcome of a single repetition of an experiment.

Definition 4.3 (Event). An event is a collection of simple events, or the outcome of multiple repetitions of an event. Events are often denoted by a capital letter.

Definition 4.4 (Mutually exclusive). Events are mutually exclusive if the occurrence of one event precludes the occurrence of another. In other words, if one event occurs, the other event cannot occur.

Definition 4.5 (Sample space). A sample space is the set of possible simple events, or all outcomes of an experiment.

4.2 Probability of Events

4.2.1 Probability of an Event

For a discrete finite sample space, the probability of a simple event is defined as the relative frequency of an outcome. Given the simple event A ,

$$\Pr(A) = \lim_{n \rightarrow \infty} \frac{I_A}{n}$$

where I_A is a function that evaluates to 1 if A occurs and 0 otherwise. The probabilities of events must satisfy the following conditions:

- $0 \leq \Pr(A) \leq 1$, where A is a simple event.
- The sum of the probabilities over the sample space is 1.

If an event A consists of a collection of simple events and each outcome is equally likely, then we can calculate the probability of an event as

$$\Pr(A) = \frac{\text{number of ways that } A \text{ can occur}}{\text{total number of outcomes}}$$

4.2.2 Probability of an Event in a Sample Space

Given the continuous sample space S , the event A can be defined as a subset of S , $A \subseteq S$. The definition of the probability of event A can be written as

$$\Pr(A) = \frac{\text{the area of region-}A}{\text{the area of region-}S}$$

As this probability is a ratio, it can be standardised so that the area of S is 1. Thus $\Pr(A)$ is the area of region- A .

4.2.3 Probability of the Complement

The complement of an event A is every event not in A , and is denoted as A^c or \bar{A} . Since the total probability for the sample space is 1, then the probability of A^c is:

$$\Pr(A^c) = 1 - \Pr(A)$$

This is true because $A \cup A^c = S$ and $\Pr(S) = 1$.

4.2.4 Probability of Subsets

If $B \subset A$, then $\Pr(B) \leq \Pr(A)$.

4.2.5 Addition Law

$$\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B).$$

and for disjoint events:

$$\Pr(A \cup B) = \Pr(A) + \Pr(B).$$

so that the intersection $\Pr(A \cap B) = 0$.

4.3 Conditional Probability

Conditional probability is the probability of an event A given another event B occurs.

If $A \cap B = \emptyset$, then $\Pr(A \cap B) = 0$. Thus if we know that B has occurred, then we know that A cannot occur:

$$\Pr(A | B) = 0.$$

Therefore if $A \cap B \neq \emptyset$, then $\Pr(A \cap B) \neq 0$. If $B \neq \emptyset$, then the conditional probability of A given B is given by:

$$\Pr(A | B) = \frac{\Pr(A \cap B)}{\Pr(B)}.$$

Note that when $\Pr(A) > \Pr(B)$, $\Pr(A | B) > \Pr(B | A)$.

4.3.1 Independence

Independence can be defined in terms of conditional probability. A and B are independent events if

$$\Pr(A | B) = \Pr(A).$$

This leads to the multiplication rule for independent events:

$$\Pr(A \cup B) = \Pr(A) \Pr(B).$$

4.4 Bayes' Rules

4.4.1 Law of Total Probability

By partitioning the sample space S into a collection of disjoint events B_1, B_2, \dots, B_n , such that $\bigcup_{i=1}^n B_i = S$, we have

$$\Pr(A) = \sum_{i=1}^n \Pr(A | B_i) \Pr(B_i)$$

where $\Pr(A | B_i) \Pr(B_i) = \Pr(A \cap B_i)$. Given the probability for A given B , the probability of the reverse direction is given by

$$\Pr(B | A) = \frac{\Pr(A | B) \Pr(B)}{\Pr(A)}$$

this is known as **Bayes' Theorem**. Using the law of total probability, we can express this as

$$\Pr(B|A) = \frac{\Pr(A|B)\Pr(B)}{\sum_{i=1}^n \Pr(A|B_i)\Pr(B_i)}.$$

5 Probability Distributions

5.1 Random Variables

A random variable is a variable whose value is the result of an experiment or random trial, where the value is not known before the trial with certainty.

5.2 Discrete Random Variables

Definition 5.1 (Discrete random variables). A discrete random variable takes on values in \mathbb{N}_0 , where the random variables arises from counting processes.

5.2.1 Probability Mass Function

The probability mass function (PMF) of a discrete random variable X is a function $p(x)$ that maps values from the sample space of X onto the interval $[0, 1]$.

$$p(x) = \Pr(X = x).$$

This function is constrained by the following properties:

- $p(x) = 0$ if $x \notin X$.
- $p(x) \in [0, 1]$ if $x \in X$.
- $\sum_{x \in X} p(x) = 1$.

5.2.2 Cumulative Mass Function

The cumulative mass function (CMF) is defined

$$F(x) = \Pr(X \leq x) = \sum_{-\infty}^x p(x)$$

and the probabilities for events can be defined using the CMF:

$$\Pr(a < X \leq b) = \sum_{x=a+1}^b p(x) = F(b) - F(a).$$

5.2.3 Expectation

The expectation (expected value) of a random variable X with a PMF is given by:

$$E(X) = \sum_{\forall x \in X} xp(x)$$

where the expectation is often denoted μ . As the expectation is a weighted average of all possible values in X , we can extend this definition to any function of X :

$$E(h(X)) = \sum_{\forall x \in X} h(x)p(x).$$

5.2.4 Median and Mode

The median m of a discrete random variable X is defined:

$$m \in X : \Pr(X \leq m) \geq \frac{1}{2} \wedge \Pr(X \geq m) \geq \frac{1}{2}.$$

Note that $\Pr(X \leq x) = \sum_{-\infty}^x p(x)$ and $\Pr(X \geq x) = \sum_x^{\infty} p(x)$.
The mode of a discrete random variable X is defined:

$$\max_{x \in X} p(x).$$

5.2.5 Variance

The variance of a random variable X is defined using the mean:

$$\begin{aligned} \text{Var}(X) &= E((X - \mu)^2) \\ &= E(X^2) - 2\mu E(X) + \mu^2 \\ &= E(X^2) - E(X)^2 \\ &= \sum_{\forall x \in X} (x - \mu)^2 p(x). \end{aligned}$$

5.3 Continuous Random Variables

Definition 5.2 (Continuous random variables). A continuous random variable takes on values in R , where values lie on a continuum.

5.3.1 Probability Density Function

The probability density function (PDF) of a continuous random variable X is a function $f(x)$ that describes the density of possible values for a continuous random variable. Note that it does not define the probability of a specific value.

For a continuous random variable X :

- $\Pr(X = x) = 0$ and $f(x) \neq \Pr(X = x)$.
- $\forall x \in X : f(x) \geq 0$.

- $\int_{-\infty}^{\infty} f(u) du = 1.$

As the probability of a single value is zero, we can instead quantify the probability of a range of values:

$$\Pr(a \leq x \leq b) = \int_a^b f(x) dx$$

where $a \leq b$.

5.3.2 Cumulative Density Function

The cumulative density function (CDF) is defined

$$F(x) = \Pr(X \leq x) = \sum_{-\infty}^x f(u) du$$

so that

$$\Pr(a \leq X \leq b) = F(b) - F(a).$$

In the continuous case, the PDF and CDF are related through the following differential equation:

$$\frac{dF}{dx} = f(x).$$

5.3.3 Median and Mode

The median m of a continuous random variable X is defined:

$$m : \int_{-\infty}^m f(u) du = \frac{1}{2}$$

and the mode is defined:

$$\max_{x \in X} f(x) \quad \text{or} \quad m : \frac{df(y)}{dy} = 0.$$

5.3.4 Expectation

The expectation of a continuous random variable X is defined as

$$\mu = E(X) = \int_{-\infty}^{\infty} xf(x) dx$$

5.3.5 Variance

The variance of a continuous random variable X is defined:

$$\text{Var}(X) = E(X^2) - E(X)^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x)^2 f(x) dx.$$

5.4 Probability Distributions

5.4.1 Bernoulli Distribution

A Bernoulli (or binary) distribution describes the probability distribution of a Boolean-valued outcome, i.e., success (1) or failure (0).

A discrete random variable X with a Bernoulli distribution is denoted

$$X \sim \text{Bernoulli}(p)$$

with

$$\begin{aligned} \Pr(X = x) &= \begin{cases} 1 - p & x = 0 \\ p & x = 1 \end{cases} \\ &= p^x (1 - p)^{1-x} \\ \Pr(X \leq x) &= \begin{cases} 0 & x < 0 \\ 1 - p & 0 \leq x < 1 \\ 1 & k \geq 1 \end{cases} \end{aligned}$$

for a probability $p \in [0, 1]$ and outcomes $x \in \{0, 1\}$. We can also summarise the following:

$$\begin{aligned} \mathbb{E}(X) &= p \\ \text{Var}(X) &= p(1 - p) \end{aligned}$$

where $(1 - p)$ is sometimes denoted as q .

5.4.2 Binomial Distribution

A binomial distribution describes the probability distribution of the number of successes for n independent trials with the same probability of success p .

A discrete random variable X with a binomial distribution is denoted

$$X \sim \text{Binomial}(n, p)$$

with

$$\begin{aligned} \Pr(X = x) &= \binom{n}{x} p^x (1 - p)^{n-x} \\ \Pr(X \leq x) &= \sum_{u=0}^x \binom{n}{u} p^u (1 - p)^{n-u} \end{aligned}$$

for number of successes $x \in \{0, 1, \dots, n\}$.

Here each individual trial is a Bernoulli trial, so that X can be written as the sum of n *independent and identically distributed* (iid) Bernoulli random variables, Y_1, Y_2, \dots, Y_n .

$$X = Y_1 + Y_2 + \dots + Y_n, \quad Y_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p) : \forall i \in \{1, 2, \dots, n\}.$$

We can then summarise the following:

$$\begin{aligned} \mathbb{E}(X) &= np \\ \text{Var}(X) &= np(1 - p) \end{aligned}$$

5.4.3 Poisson Distribution

A Poisson distribution describes the probability distribution of the number of events N which occur over a fixed interval of time λ .

A discrete random variable N with a Poisson distribution is denoted

$$N \sim \text{Poisson}(\lambda)$$

with

$$\begin{aligned}\Pr(N = n) &= \frac{\lambda^n e^{-\lambda}}{n!} \\ \Pr(N \leq n) &= e^{-\lambda} \sum_{u=0}^n \frac{\lambda^u}{u!}\end{aligned}$$

for number of events $n \geq 0$. We can also summarise the following:

$$\begin{aligned}\mathbb{E}(N) &= \lambda \\ \text{Var}(N) &= \lambda\end{aligned}$$

The Poisson PMF can be defined in terms of the Binomial PMF as $n \rightarrow \infty$ and $p \rightarrow 0$. Let $\lambda = np$, then

$$\begin{aligned}p(x) &= \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x} \\ &= \frac{n!}{x!(n-x)!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^{n-x} \\ &= \frac{\lambda^x}{x!} \frac{n!}{(n-x)!} \frac{1}{n^x} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x}\end{aligned}$$

The limit of $\frac{n!}{(n-x)!} \frac{1}{n^x}$ is 1:

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{n!}{(n-x)!} \frac{1}{n^x} &= \lim_{n \rightarrow \infty} \frac{n(n-1)(n-2) \cdots (n-x+1)}{n^x} \\ &= \lim_{n \rightarrow \infty} \left(\frac{n}{n}\right) \left(\frac{n-1}{n}\right) \left(\frac{n-2}{n}\right) \cdots \left(\frac{n-x+1}{n}\right) \\ &= 1\end{aligned}$$

The term $\left(1 - \frac{\lambda}{n}\right)^n$ approaches $e^{-\lambda}$, using the substitution $u = -\frac{n}{\lambda}$:

$$\begin{aligned}\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^n &= \lim_{n \rightarrow \infty} \left(1 + \frac{1}{u}\right)^{u(-\lambda)} \\ &= e^{-\lambda}\end{aligned}$$

Finally, the remaining term also evaluates to 1:

$$\lim_{n \rightarrow \infty} \left(1 - \frac{\lambda}{n}\right)^{-x} = 1$$

Therefore by gathering the above equations, we can write the Poisson PMF as:

$$\begin{aligned} p(x) &= \lim_{n \rightarrow \infty} \frac{\lambda^x}{x!} \frac{n!}{(n-x)!} \frac{1}{n^x} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x} \\ &= \frac{\lambda^x e^{-\lambda}}{x!} \end{aligned}$$

5.4.4 Uniform Distribution

A continuous uniform distribution describes the probability distribution of an outcome within some interval, where the probability of an outcome in one interval is the same as all other intervals of the same length.

A continuous random variable X with a continuous uniform distribution is denoted

$$X \sim \text{Uniform}(a, b)$$

with

$$\begin{aligned} f(x) &= \frac{1}{b-a} \\ F(x) &= \frac{x-a}{b-a} \end{aligned}$$

for outcomes $a < x < b$. We can also summarise the following:

$$\begin{aligned} E(X) &= \frac{a+b}{2} \\ \text{Var}(X) &= \frac{(b-a)^2}{12} \\ m &= \frac{a+b}{2} \end{aligned}$$

5.4.5 Exponential Distribution

An exponential distribution describes the probability distribution of the time between events with rate η .

A continuous random variable T with an exponential distribution is denoted

$$T \sim \text{Exp}(\eta)$$

with

$$\begin{aligned} f(t) &= \eta e^{-\eta t} \\ F(t) &= 1 - e^{-\eta t} \end{aligned}$$

for time $t > 0$. We can also summarise the following:

$$\begin{aligned} E(T) &= \frac{1}{\eta} \\ \text{Var}(T) &= \frac{1}{\eta^2} \\ m &= \frac{\ln(2)}{\eta} \end{aligned}$$

Proof. By considering an event taking longer than t seconds, we can represent this as nothing happening over the interval $[0, t]$. Using $T \sim \text{Exp}(\eta)$ and $N \sim \text{Poisson}(\eta t)$, we have

$$\Pr(T > t) = \Pr(N = 0) = e^{-\eta t}$$

where $\lambda = \eta t$. The CDF for the exponential distribution is then

$$\begin{aligned}\Pr(T < t) &= 1 - \Pr(T > t) \\ &= 1 - e^{-\eta t}.\end{aligned}$$

□

5.4.6 Memoryless Property

In an exponential distribution with $T \sim \text{Exp}(\eta)$, the distribution of the waiting time $t + s$ until a certain event does not depend on how much time t has already passed.

$$\Pr(T > s + t \mid T > t) = \Pr(T > s).$$

The same property also applies in an Geometric distribution with $N \sim \text{Geometric}(p)$.

5.4.7 Normal Distribution

The normal distribution is used to represent many random situations, in particular, measurements and their errors. This distribution arises in many statistical problems and can be used to approximate other distributions under certain conditions.

A continuous random variable X with a normal distribution is denoted

$$X \sim \text{N}(\mu, \sigma^2)$$

with

$$\begin{aligned}f(t) &= \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \\ F(t) &= \frac{1}{2} \left(1 + \text{erf}\left(\frac{x-\mu}{\sigma\sqrt{2}}\right) \right)\end{aligned}$$

for $x \in \mathbb{R}$ where $\text{erf}(z) = \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt$ is the error function. We can also summarise the following:

$$\begin{aligned}\text{E}(X) &= \mu \\ \text{Var}(X) &= \sigma^2\end{aligned}$$

Given the complexity of the analytic expressions for the PDF and CDF of the normal distribution, we often use software to numerically determine probabilities associated with normal distributions.

5.5 Standard Normal Distribution

Given $X \sim \text{N}(\mu, \sigma^2)$, consider the transformation

$$Z = \frac{X - \mu}{\sigma}$$

so that $Z \sim \text{N}(0, 1)$. This distribution is called the standard normal distribution. This allows us to deal with the standard normal distribution regardless of μ and σ .

6 Sampling

This section explores the ideas of sampling and inference.

6.1 Observational and Experimental Studies

6.1.1 Observational Studies

If we are collecting or sampling data that already exists, i.e., we have no control over how we created the data, then we are conducting an observational study.

6.1.2 Experimental Studies

If data was generated in a controlled experimental environment, then the data is the result of an experimental study.

6.2 Sampling

As it is impossible to collect measurements from an entire population, we must rely on samples and sample statistics to make inferences about the population. There are several methods for this, depending on the situation.

6.2.1 Simple Random Sampling

In simple random sampling, a random subset of size n is selected from a population of size N . Simple random sampling is used in **observational studies** as data already exists. There are some caveats to this method that may lead to errors in the conclusions reached:

- Non-response bias: some members may not respond to the survey.
- Undercoverage bias: the survey may not apply to all members of the selection.
- Wording bias: the wording of the survey may lead to a biased response.

6.3 Stratified Random Sampling

Stratified random sampling is a method of sampling that divides the population into non-overlapping strata and draws random samples from each stratum.

- **Pre-stratification** is when the strata are defined before the sampling process.
- **Post-stratification** is when the strata are defined after the sampling process.

6.4 Cluster Sampling

Cluster sampling is used when there are limited resources or a lack of information about individuals in the population. It is also useful when members in each cluster are similar.

6.5 Non-random Sampling Methods

- Sequential sampling: samples are taken in a sequential manner.
- Convenience sampling: samples are self-selected from the most convenient source.
- Snowball sampling is like convenience sampling but participants are asked to recruit others.
- Quota sampling: samples are selected to balance a particular demographic.

7 Sampling Distributions

A sampling distribution is the probability distribution of a sample statistic.

7.1 Central Limit Theorem

For a sample of size n from any random probability distribution with expected value μ and variance σ^2 ,

$$\frac{\sqrt{n}(\bar{x} - \mu)}{\sigma} \xrightarrow{p} N(0, 1)$$

meaning that increasing the sample size will lead to a more normal distribution. In this case, a sample size of $n = 30$ is sufficient to approximate a normal distribution.

7.2 Standard Error

The standard error of a sample statistic is the standard deviation of the sampling distribution.

$$SE(\bar{x}) = \frac{\sigma^2}{n}$$

7.3 Sample Proportion

Sometimes we are interested in estimating the population proportion from the sample proportion. For a sample of size n let x be the number of members with a particular characteristic. The sample estimate of the population proportion p is

$$\hat{p} = \frac{x}{n}.$$

By assuming that the sample statistic x follows a binomial distribution with probability p and size n , then $E(x) = np$ and $\text{Var}(x) = np(1-p)$. Therefore the expectation is

$$E(\hat{p}) = p$$

and the standard error is

$$SE(\hat{p}) = \sqrt{\frac{p(1-p)}{n}}.$$

For the above to apply, we must assume that the sample proportion and size are sufficiently large. In general, if $np > 5$ and $n(1-p) > 5$, then we can assume that the sampling distribution of \hat{p} is approximately normal.

7.4 Assessing Normality

- Histograms: if the data is approximately normal, then the histogram will be approximately symmetric and unimodal.
- Boxplots: boxplots can be useful for showing outliers and skewness. Extreme clusters of an excessive number of outliers can be evidence of non-normality.
- Normal probability plots (q - q plots): these plots are constructed by plotting the sorted data values against their Z -scores. If the data is approximately normal, then the points will lie approximately on a straight line.