

# Statistical Inference and Modelling Summary

Daniel Weroniski Falcó

October 2024

## Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>   | <b>3</b>  |
| 1.1      | Classical vs Fisherian Inference . . . . .                    | 3         |
| 1.2      | Key Concepts in Inference . . . . .                           | 3         |
| 1.3      | When to Use Different Methods . . . . .                       | 3         |
| 1.4      | Useful R Functions for Inference . . . . .                    | 3         |
| <b>2</b> | <b>Exploratory Data Analysis</b>                              | <b>4</b>  |
| 2.1      | Numerical Summaries . . . . .                                 | 4         |
| 2.2      | Graphical Summaries . . . . .                                 | 4         |
| 2.3      | Useful R Functions for EDA . . . . .                          | 4         |
| 2.4      | Example R Code for EDA . . . . .                              | 4         |
| 2.5      | Key Points . . . . .  | 5         |
| <b>3</b> | <b>Hypothesis Testing</b>                                     | <b>6</b>  |
| 3.1      | Key Concepts in Hypothesis Testing . . . . .                  | 6         |
| 3.2      | Types of Hypothesis Tests . . . . .                           | 6         |
| 3.3      | Useful R Functions for Hypothesis Testing . . . . .           | 6         |
| 3.4      | Example R Code for Hypothesis Testing . . . . .               | 7         |
| 3.5      | General Procedure for Hypothesis Testing . . . . .            | 7         |
| 3.6      | Choosing the Appropriate Hypothesis Test . . . . .            | 8         |
| <b>4</b> | <b>Confidence Intervals</b>                                   | <b>9</b>  |
| 4.1      | Key Concepts in Confidence Intervals . . . . .                | 9         |
| 4.2      | Confidence Interval for the Mean (Known Variance) . . . . .   | 9         |
| 4.3      | Confidence Interval for the Mean (Unknown Variance) . . . . . | 9         |
| 4.4      | Confidence Interval for Proportions . . . . .                 | 9         |
| 4.5      | Useful R Functions for Confidence Intervals . . . . .         | 10        |
| 4.6      | Example R Code for Confidence Intervals . . . . .             | 10        |
| 4.7      | Choosing the Appropriate Confidence Interval . . . . .        | 10        |
| 4.8      | Key Points . . . . .  | 10        |
| <b>5</b> | <b>Distributions</b>  | <b>11</b> |
| 5.1      | Key Distributions and Their Properties . . . . .              | 11        |
| 5.2      | Useful R Functions for Distributions . . . . .                | 11        |
| 5.3      | Example R Code for Distributions . . . . .                    | 12        |
| 5.4      | When to Use Each Distribution . . . . .                       | 12        |

|          |  |           |
|----------|--|-----------|
| <b>6</b> | <b>Regression and Modelling</b>                            | <b>13</b> |
| 6.1      | Linear Regression . . . . .                                | 13        |
| 6.1.1    | Multiple Linear Regression . . . . .                       | 13        |
| 6.1.2    | Assumptions of Linear Regression . . . . .                 | 13        |
| 6.2      | Logistic Regression . . . . .                              | 13        |
| 6.3      | Model Evaluation Metrics . . . . .                         | 14        |
| 6.4      | Useful R Functions for Regression and Modelling . . . . .  | 14        |
| 6.5      | Example R Code for Regression Analysis . . . . .           | 14        |
| 6.6      | Choosing Between Models . . . . .                          | 15        |
| 6.7      | When to Use Different Regression Models . . . . .          | 15        |
| 6.8      | Key Points . . . . .                                       | 15        |
| <b>7</b> | <b>Advanced Topics</b>                                     | <b>16</b> |
| 7.1      | Central Limit Theorem (CLT) . . . . .                      | 16        |
| 7.2      | Law of Large Numbers (LLN) . . . . .                       | 16        |
| 7.3      | Non-Parametric Tests . . . . .                             | 16        |
| 7.4      | Bootstrap Methods . . . . .                                | 17        |
| 7.5      | Choosing Between Advanced Techniques . . . . .             | 17        |
| 7.6      | Key Points . . . . .                                       | 17        |
| <b>8</b> | <b>Useful R Functions and Syntax</b>                       | <b>19</b> |
| 8.1      | Data Handling and Manipulation . . . . .                   | 19        |
| 8.2      | Basic Plotting Functions . . . . .                         | 19        |
| 8.3      | Advanced Visualisation with <code>ggplot2</code> . . . . . | 19        |
| 8.4      | Statistical Modelling Functions . . . . .                  | 19        |
| 8.5      | Cheat Sheet Table: Common R Functions . . . . .            | 20        |
| 8.6      | R Markdown for Reproducible Reports . . . . .              | 20        |
| 8.7      | Additional Resources and Tips . . . . .                    | 20        |
| <b>9</b> | <b>Quick Overview of Statistical Tests</b>                 | <b>21</b> |
| 9.1      | Choosing the Appropriate Test . . . . .                    | 21        |
| 9.2      | Overview of Hypothesis Tests . . . . .                     | 21        |
| 9.3      | Example R Code for Various Tests . . . . .                 | 23        |

# 1 Introduction

This cheat sheet provides a concise summary of key concepts in statistical inference and modeling. It is structured to help navigate through definitions, usages, and code snippets for statistical tools commonly used in data science.

## 1.1 Classical vs Fisherian Inference

Statistical inference is broadly classified into two paradigms:

- **Classical Inference (Frequentist Approach):** Focuses on the long-run frequency of events. Relies heavily on hypothesis testing, confidence intervals, and p-values.
- **Fisherian Inference:** An approach that uses likelihood to measure support for different hypotheses. Unlike the classical approach, Fisher's inference does not require specification of alternative hypotheses.

## 1.2 Key Concepts in Inference

Statistical inference involves making decisions or drawing conclusions about a population based on a sample. The two primary goals are:

1. **Estimation:** Determining the approximate value of a population parameter.
2. **Hypothesis Testing:** Assessing evidence provided by the data against a specified hypothesis.

## 1.3 When to Use Different Methods

- **Parametric Tests** (e.g., t-tests, ANOVA): Assume underlying distribution of the data.
- **Non-Parametric Tests** (e.g., Mann-Whitney U test, Kolmogorov-Smirnov test, Kruskal-Wallis test): Do not assume any specific distribution.

## 1.4 Useful R Functions for Inference

- `t.test()` — Performs one and two-sample t-tests.
- `chisq.test()` — Performs Chi-squared test for independence.
- `prop.test()` — Tests for the equality of proportions.

## 2 Exploratory Data Analysis

Exploratory Data Analysis (EDA) is a crucial step in understanding and summarizing the main characteristics of a dataset. It involves the use of various statistical tools and visualization techniques to identify patterns, detect outliers, and gain insights before conducting formal modelling or hypothesis testing.

### 2.1 Numerical Summaries

**Numerical summaries** provide quantitative insights into the dataset using measures of central tendency and variability. Key statistics include:

- **Mean** ( $\mu$ ): The average of all observations; the sum of all possible values weighted by their probability:  $\mu = \langle X \rangle = \sum_{i=1}^n x_i p(x_i)$
- **Median** ( $Q_2$ ): The middle value of the ordered dataset, less sensitive to outliers.
- **Mode**: The most frequently occurring value(s) in the dataset.
- **Variance** ( $\sigma^2$ ): Measure of how far observations spread out from the mean. Measured as the expectation value of the square deviation from the mean:  $\sigma^2 = \langle (X - \mu)^2 \rangle = \langle X^2 \rangle - \langle X \rangle^2$
- **Standard Deviation** ( $\sigma$ ): The square root of the variance.
- **Interquartile Range** (IQR): Difference between the 75th percentile ( $Q_3$ ) and 25th percentile ( $Q_1$ ).

### 2.2 Graphical Summaries

Graphical representations help to visualize the distribution and relationships between variables. Common plots include:

- **Histograms**: Show the frequency distribution of a continuous variable.
- **Box Plots**: Display the spread of the data along with outliers.
- **Scatter Plots**: Used to visualize the relationship between two continuous variables.
- **Bar Charts**: Represent categorical data using bars of different heights.

### 2.3 Useful R Functions for EDA

- `summary()` — Provides a numerical summary of a dataset.
- `boxplot()` — Plots a box-and-whisker plot.
- `hist()` — Creates a histogram for a given variable.
- `plot()` — Generic plotting function, often used for scatter plots.

### 2.4 Example R Code for EDA

Below is an example of performing basic EDA on a sample dataset in R:

```
1 # Load necessary library
2 library(ggplot2)
3
4 # Example dataset: mtcars
5 data("mtcars")
6
```

```

7 # Numerical summary
8 summary(mtcars)
9
10 # Histogram of Miles per Gallon (mpg)
11 hist(mtcars$mpg, main="Histogram of MPG", xlab="Miles Per Gallon",
12      col="lightblue", breaks=10)
13
14 # Boxplot for distribution of mpg by number of cylinders
15 boxplot(mpg ~ cyl, data=mtcars, main="Boxplot of MPG by Cylinder",
16        xlab="Number of Cylinders", ylab="Miles Per Gallon", col="orange")
17
18 # Scatter plot of Horsepower vs MPG
19 plot(mtcars$hp, mtcars$mpg, main="Scatter Plot of Horsepower vs MPG",
20      xlab="Horsepower", ylab="Miles Per Gallon", col="blue", pch=19)

```

Listing 1: Basic EDA in R

## 2.5 Key Points

- EDA helps to identify anomalies and understand the structure of the data.
- Both numerical and graphical summaries are essential in the preliminary stages of data analysis.
- Visualizations are often more intuitive and revealing compared to raw statistical measures.

## 3 Hypothesis Testing

Hypothesis testing is a statistical method used to make decisions or draw conclusions about a population based on sample data. It involves formulating two competing hypotheses and using sample data to determine which hypothesis is more plausible.

### 3.1 Key Concepts in Hypothesis Testing

- **Null Hypothesis ( $H_0$ ):** The hypothesis that there is no effect or no difference. This is the hypothesis that we aim to test against.
- **Alternative Hypothesis ( $H_1$ ):** The hypothesis that there is a significant effect or difference.
- **Significance Level ( $\alpha$ ):** The probability of rejecting the null hypothesis when it is actually true. Common values are 0.05, 0.01, or 0.1.
- **p-value:** The probability of obtaining results as extreme as those observed, assuming  $H_0$  is true. A small p-value (less than  $\alpha$ ) indicates strong evidence against  $H_0$ .
- **Type I Error:** Incorrectly rejecting  $H_0$  when it is true (false positive).
- **Type II Error:** Failing to reject  $H_0$  when  $H_1$  is true (false negative).

### 3.2 Types of Hypothesis Tests

The choice of test depends on the nature of the data and the research question. Some common hypothesis tests include:

- **t-tests** — Used to compare means.
  - **One-sample t-test:** Compares the mean of a single group to a known value.
  - **Two-sample t-test:** Compares means between two independent groups.
  - **Paired t-test:** Compares means from the same group at different times.
- **ANOVA (Analysis of Variance)** — Compares means among three or more groups.
- **Chi-squared Test** — Assesses relationships between categorical variables.
- **Fisher's Exact Test** — An alternative to Chi-squared when sample sizes are small.

### 3.3 Useful R Functions for Hypothesis Testing

- `t.test()` — Performs one and two-sample t-tests, as well as paired t-tests.
- `aov()` — Fits an ANOVA model.
- `chisq.test()` — Performs the Chi-squared test for independence.
- `fisher.test()` — Conducts Fisher's Exact Test for count data.
- `wilcox.test()` — Performs non-parametric tests (e.g., Mann-Whitney U test).

### 3.4 Example R Code for Hypothesis Testing

The following example shows how to perform various hypothesis tests in R:

```
1 # Load necessary data
2 data(mtcars)
3
4 # One-sample t-test
5 t.test(mtcars$mpg, mu = 20) # Test if mean mpg is different from 20
6
7 # Two-sample t-test
8 t.test(mpg ~ cyl, data = subset(mtcars, cyl %in% c(4, 6))) # Test if
   mean mpg differs by number of cylinders (4 vs. 6)
9
10 # Paired t-test
11 pre_weight <- c(70, 75, 80, 78, 72)
12 post_weight <- c(68, 74, 79, 76, 71)
13 t.test(pre_weight, post_weight, paired = TRUE) # Compare pre and post
   weights
14
15 # Chi-squared Test for independence
16 chisq.test(mtcars$cyl, mtcars$gear)
17
18 # ANOVA
19 fit <- aov(mpg ~ factor(cyl), data = mtcars) # Test if mpg differs
   across cylinders
20 summary(fit)
```

Listing 2: Hypothesis Testing in R

### 3.5 General Procedure for Hypothesis Testing

The standard approach for hypothesis testing involves the following steps:

**1. State the Hypotheses:**

- Null hypothesis ( $H_0$ ): Example: There is no difference in mean mileage between 4-cylinder and 6-cylinder cars.
- Alternative hypothesis ( $H_1$ ): Example: The mean mileage is different between the two groups.

**2. Set the Significance Level ( $\alpha$ ):** Common choices are 0.05 or 0.01.

**3. Calculate the Test Statistic:** Use the appropriate formula depending on the test being conducted.

**4. Find the p-value and Make a Decision:**

- If  $p \leq \alpha$ : Reject  $H_0$ , and conclude that there is sufficient evidence to support  $H_1$ .
- If  $p > \alpha$ : Fail to reject  $H_0$ , and conclude that there is not enough evidence to support  $H_1$ .

### 3.6 Choosing the Appropriate Hypothesis Test

| Scenario   | Test to Use       | R Function                                    |
|--|-------------------|---|
| Comparing the mean of a sample to a known value          | One-sample t-test | <code>t.test(x, mu = value)</code>            |
| Comparing means between two independent groups           | Two-sample t-test | <code>t.test(x group)</code>                  |
| Comparing means within the same group at different times | Paired t-test     | <code>t.test(pre, post, paired = TRUE)</code> |
| Testing for relationship between categorical variables   | Chi-squared Test  | <code>chisq.test(table)</code>                |
| Comparing means across multiple groups                   | ANOVA             | <code>aov()</code>                            |



## 4 Confidence Intervals

Confidence intervals provide a range of plausible values for an unknown population parameter based on sample data. They are used to quantify the uncertainty around a point estimate and are typically expressed as a percentage (e.g., 95% confidence interval).

### 4.1 Key Concepts in Confidence Intervals

- **Point Estimate:** A single value estimate of a population parameter (e.g., sample mean).
- **Confidence Level:** The proportion of times the interval would contain the true parameter if we were to repeat the sampling process many times (e.g., 95%).
- **Margin of Error:** The maximum expected difference between the point estimate and the true population parameter. Influenced by sample size and variability.
- **Confidence Interval Formula:**

$$\text{CI} = \text{Point Estimate} \pm \text{Critical Value} \times \text{Standard Error}$$

### 4.2 Confidence Interval for the Mean (Known Variance)

When the population variance  $\sigma^2$  is known, the confidence interval for the mean  $\mu$  is calculated as:

$$\text{CI} = \bar{x} \pm z_{\alpha/2} \times \frac{\sigma}{\sqrt{n}}$$

Where:

- $\bar{x}$  — Sample mean.
- $z_{\alpha/2}$  — Critical value from the standard normal distribution for the given confidence level.
- $\sigma$  — Population standard deviation.
- $n$  — Sample size.

### 4.3 Confidence Interval for the Mean (Unknown Variance)

When the population variance is unknown, use the sample standard deviation  $s$  and the t-distribution:

$$\text{CI} = \bar{x} \pm t_{\alpha/2, n-1} \times \frac{s}{\sqrt{n}}$$

Where  $t_{\alpha/2, n-1}$  is the critical value from the t-distribution with  $n - 1$  degrees of freedom.

### 4.4 Confidence Interval for Proportions

For categorical data, the confidence interval for a proportion  $p$  is calculated using:

$$\text{CI} = \hat{p} \pm z_{\alpha/2} \times \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

Where:

- $\hat{p}$  — Sample proportion.
- $n$  — Sample size.
- $z_{\alpha/2}$  — Critical value from the standard normal distribution.

## 4.5 Useful R Functions for Confidence Intervals

- `t.test()` — Provides confidence intervals for one-sample, two-sample, and paired t-tests.
- `prop.test()` — Computes confidence intervals for proportions.
- `confint()` — Calculates confidence intervals for model parameters.

## 4.6 Example R Code for Confidence Intervals

Here's how to compute various types of confidence intervals in R:

```
1 # One-sample t-test for confidence interval
2 x <- c(12.9, 14.2, 13.7, 15.3, 14.8, 13.6, 14.0)
3 t.test(x) # 95% CI for the mean
4
5 # Confidence interval for a proportion
6 prop.test(x = 34, n = 230, conf.level = 0.95) # 95% CI for proportion
7
8 # Confidence interval for a difference in means
9 group1 <- c(15, 17, 16, 14, 18)
10 group2 <- c(12, 14, 13, 15, 11)
11 t.test(group1, group2) # 95% CI for the difference in means
12
13 # Confidence intervals for regression coefficients
14 model <- lm(mpg ~ hp + wt, data = mtcars)
15 confint(model, level = 0.95) # 95% CI for model coefficients
```

Listing 3: Confidence Interval Calculations in R

## 4.7 Choosing the Appropriate Confidence Interval

| Scenario                              | Parameter of Interest | Distribution          | R Function                                 |
|---------------------------------------|-----------------------|-----------------------|--|
| Population mean with known variance   | $\mu$                 | Normal distribution   | Manual calculation or <code>qnorm()</code> |
| Population mean with unknown variance | $\mu$                 | t-distribution        | <code>t.test()</code>                      |
| Population proportion                 | $p$                   | Binomial distribution | <code>prop.test()</code>                   |
| Difference in two population means    | $\mu_1 - \mu_2$       | t-distribution        | <code>t.test()</code>                      |
| Regression coefficients               | $\beta$               | Normal approximation  | <code>confint()</code>                     |

## 4.8 Key Points

- Higher confidence levels (e.g., 99%) produce wider intervals compared to lower confidence levels (e.g., 90%).
- Confidence intervals are not probabilities; they describe the range of plausible values for the population parameter.
- If a confidence interval for a mean difference contains zero, it suggests no significant difference.

## 5 Distributions

Understanding different probability distributions is crucial in statistical modelling and inference. Each distribution has unique properties and applications, and choosing the correct one for your data is essential.

### 5.1 Key Distributions and Their Properties

- **Normal Distribution:**

- Symmetrical, bell-shaped distribution defined by its mean  $\mu$  and variance  $\sigma^2$ .
- Widely used in inferential statistics due to the Central Limit Theorem.
- **R Functions:** `dnorm()`, `pnorm()`, `rnorm()` (see section 5.2).

- **Binomial Distribution:**

- Models the number of successes in a fixed number of independent Bernoulli trials.
- Defined by parameters  $n$  (number of trials) and  $p$  (probability of success).
- **R Functions:** `dbinom()`, `pbinom()`, `rbinom()`.

- **Poisson Distribution:**

- Models the count of events occurring in a fixed interval of time or space.
- Parameterized by  $\lambda$  (rate of occurrence).
- **R Functions:** `dpois()`, `ppois()`, `rpois()`.

- **Exponential Distribution:**

- Models the time between successive events in a Poisson process.
- Defined by the rate parameter  $\lambda$ .
- **R Functions:** `dexp()`, `pexp()`, `rexp()`.

- **Chi-squared Distribution:**

- Distribution of the sum of squared standard normal variables.
- Commonly used in hypothesis testing and constructing confidence intervals.
- **R Functions:** `dchisq()`, `pchisq()`, `rchisq()`.

- **Student's t Distribution:**

- Similar to the normal distribution but with heavier tails.
- Useful for small sample sizes or when population variance is unknown.
- **R Functions:** `dt()`, `pt()`, `rt()`.

### 5.2 Useful R Functions for Distributions

- `dnorm()`, `dpois()`, `dbinom()` — Probability density functions.
- `pnorm()`, `ppois()`, `pbinom()` — Cumulative distribution functions.
- `rnorm()`, `rpois()`, `rbinom()` — Random number generation.
- `qnorm()`, `qpois()`, `qbinom()` — Quantile functions.

### 5.3 Example R Code for Distributions

Here's how to generate and visualize different distributions in R:

```
1 # Set up parameters
2 mu <- 0
3 sigma <- 1
4 lambda <- 2
5 n <- 10
6 p <- 0.3
7
8 # Generate and visualize a normal distribution
9 x <- seq(-3, 3, by = 0.1)
10 y <- dnorm(x, mean = mu, sd = sigma)
11 plot(x, y, type = "l", main = "Normal Distribution", xlab = "X", ylab = "Density")
12
13 # Generate and plot a binomial distribution
14 binom_data <- dbinom(0:n, size = n, prob = p)
15 barplot(binom_data, main = "Binomial Distribution (n=10, p=0.3)", xlab = "Number of Successes")
16
17 # Generate and plot a Poisson distribution
18 pois_data <- dpois(0:20, lambda = lambda)
19 barplot(pois_data, main = "Poisson Distribution (lambda=2)", xlab = "Number of Events")
20
21 # Generate and plot an exponential distribution
22 exp_data <- rexp(1000, rate = lambda)
23 hist(exp_data, main = "Exponential Distribution (rate=2)", xlab = "Time", breaks = 20)
```

Listing 4: R Code for Generating and Visualizing Distributions

### 5.4 When to Use Each Distribution

| Distribution | Type of Data | Common Usage   | R Functions                                   |
|--------------|--------------|--|---|
| Normal       | Continuous   | Modeling measurement errors, heights, weights        | <code>rnorm()</code> , <code>dnorm()</code>   |
| Binomial     | Discrete     | Modeling binary outcomes (e.g., coin flips)          | <code>rbinom()</code> , <code>dbinom()</code> |
| Poisson      | Discrete     | Modeling count data (e.g., number of calls per hour) | <code>rpois()</code> , <code>dpois()</code>   |
| Exponential  | Continuous   | Modeling time to next event (e.g., failure times)    | <code>rexp()</code> , <code>dexp()</code>     |
| Chi-squared  | Continuous   | Goodness-of-fit tests                                | <code>rchisq()</code> , <code>dchisq()</code> |
| Student's t  | Continuous   | Small sample data, unknown variance                  | <code>rt()</code> , <code>dt()</code>         |

Some statistical tests rely on specific distributions; always verify assumptions before proceeding.

## 6 Regression and Modelling

Regression analysis is a powerful statistical tool used to understand the relationship between one or more independent variables and a dependent variable. It is widely used for prediction, inference, and uncovering causal relationships.

### 6.1 Linear Regression

Linear regression models the relationship between a dependent variable  $y$  and one or more independent variables  $x_1, x_2, \dots, x_n$  using a linear function.

**Simple Linear Regression Model:**

$$y = \beta_0 + \beta_1 x + \epsilon$$

Where:

- $\beta_0$  — Intercept (the value of  $y$  when  $x = 0$ ).
- $\beta_1$  — Slope (change in  $y$  for a one-unit change in  $x$ ).
- $\epsilon$  — Random error term, assumed to be normally distributed with mean 0.

#### 6.1.1 Multiple Linear Regression

Multiple linear regression models the relationship between a dependent variable and multiple independent variables:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

Where  $\beta_1, \beta_2, \dots, \beta_n$  are the coefficients for each independent variable.

#### 6.1.2 Assumptions of Linear Regression

- **Linearity:** The relationship between the independent and dependent variables is linear.
- **Independence:** Observations are independent.
- **Homoscedasticity:** Constant variance of the error terms across all levels of independent variables.
- **Normality:** The residuals (errors) follow a normal distribution.

### 6.2 Logistic Regression

Logistic regression is used when the dependent variable is categorical, often binary (e.g., success/failure). Instead of modelling the dependent variable directly, logistic regression models the log odds of the probability of an event occurring:

$$\log \left( \frac{p}{1-p} \right) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

Where:

- $p$  — Probability of success (e.g.,  $p(\text{event} = 1)$ ).
- $\beta_0, \beta_1, \dots, \beta_n$  — Coefficients estimated using maximum likelihood.

## 6.3 Model Evaluation Metrics

Once a model is fitted, it is essential to evaluate its performance using appropriate metrics:

- **R-squared** — Proportion of variance explained by the model. Higher values indicate better fit.
- **Adjusted R-squared** — Adjusted version of R-squared that penalises adding irrelevant predictors.
- **Mean Squared Error (MSE)** — Average of the squared differences between observed and predicted values.
- **Akaike Information Criterion (AIC)** — Measures the goodness of fit while penalising for the number of predictors.

## 6.4 Useful R Functions for Regression and Modelling

- `lm()` — Fits a linear regression model.
- `summary()` — Summarises the results of a linear or logistic model.
- `glm()` — Fits a generalised linear model, including logistic regression.
- `predict()` — Makes predictions using a fitted model.
- `anova()` — Compares models using ANOVA.

## 6.5 Example R Code for Regression Analysis

Here is an example of fitting linear and logistic regression models in R:

```
1 # Load the dataset
2 data(mtcars)
3
4 # Linear Regression: Predicting MPG from Horsepower and Weight
5 linear_model <- lm(mpg ~ hp + wt, data = mtcars)
6 summary(linear_model)
7
8 # Visualising the fit
9 plot(mtcars$hp, mtcars$mpg, main = "MPG vs Horsepower", xlab = "
   Horsepower", ylab = "MPG")
10 abline(linear_model, col = "red")
11
12 # Logistic Regression: Predicting Automatic vs Manual Transmission
13 mtcars$am <- factor(mtcars$am, labels = c("Automatic", "Manual"))
14 logistic_model <- glm(am ~ hp + wt, data = mtcars, family = "binomial
   ")
15 summary(logistic_model)
16
17 # Predict probabilities
18 predict(logistic_model, type = "response")
```

Listing 5: Regression Analysis in R

## 6.6 Choosing Between Models

When building regression models, consider the following:

- **Linear vs Logistic Regression:**
  - Use linear regression when the dependent variable is continuous.
  - Use logistic regression for binary outcomes.
- **Regularisation (Ridge/Lasso):**
  - Use regularisation techniques when dealing with high-dimensional data to prevent overfitting.
  - **R Function:** `glmnet()` for regularised models.

## 6.7 When to Use Different Regression Models

| Model Type             | Dependent Variable Type | Purpose                                       | R Function                            |
|------------------------|-------------------------|---|---------------------------------------|
| Linear Regression      | Continuous              | Predict a continuous outcome                  | <code>lm()</code>                     |
| Logistic Regression    | Binary/Categorical      | Classify into two categories                  | <code>glm(family = "binomial")</code> |
| Ridge/Lasso Regression | Continuous              | Prevent overfitting in high-dimensional data  | <code>glmnet()</code>                 |
| Poisson Regression     | Count Data              | Model count outcomes (e.g., number of events) | <code>glm(family = "poisson")</code>  |

## 6.8 Key Points

- Always check the assumptions of your regression model (e.g., linearity, homoscedasticity).
- Evaluate model performance using appropriate metrics (e.g., R-squared, AIC).
- Visualise residuals to check for patterns and validate model assumptions.

## 7 Advanced Topics

This section covers some of the more complex statistical concepts and methods that go beyond basic hypothesis testing and regression. Understanding these topics is essential for tackling more sophisticated analytical problems and gaining deeper insights from data.

### 7.1 Central Limit Theorem (CLT)

The Central Limit Theorem is a fundamental concept in statistics that states that the distribution of the sample mean will approximate a normal distribution, regardless of the original distribution of the population, as the sample size becomes large.

**Mathematical Statement:**

If  $X_1, X_2, \dots, X_n$  are independent and identically distributed random variables with mean  $\mu$  and variance  $\sigma^2$ , the sample mean  $\bar{X}$  is approximately normally distributed:

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

**Key Points:**

- The larger the sample size, the better the approximation to the normal distribution.
- This property justifies the use of normal-based methods (e.g., t-tests) even when the underlying data is not normally distributed.

### 7.2 Law of Large Numbers (LLN)

The Law of Large Numbers states that as the sample size increases, the sample mean  $\bar{X}$  converges to the true population mean  $\mu$ .

**Types of LLN:**

- **Weak Law:** The sample mean converges in probability to the population mean.
- **Strong Law:** The sample mean converges almost surely to the population mean.

**Implications:**

- Ensures that larger samples provide more reliable estimates of population parameters.
- The variability of the sample mean decreases as the sample size increases.

### 7.3 Non-Parametric Tests

Non-parametric tests do not assume a specific distribution for the data, making them useful for small samples or data that do not meet parametric assumptions.

- **Mann-Whitney U Test:**

- Compares the ranks of two independent samples.
- Analogous to the two-sample t-test but does not require normality.
- **R Function:** `wilcox.test()`.

- **Kruskal-Wallis Test:**

- Extends the Mann-Whitney U test to compare more than two groups.
- Equivalent to one-way ANOVA for ranked data.
- **R Function:** `kruskal.test()`.



## 7.4 Bootstrap Methods

Bootstrap methods involve resampling with replacement from a dataset to estimate the distribution of a statistic. They are useful for constructing confidence intervals and assessing the stability of estimates without relying on parametric assumptions.

### Procedure:

1. Draw  $B$  bootstrap samples from the original dataset.
2. Calculate the statistic of interest (e.g., mean, median) for each bootstrap sample.
3. Use the distribution of these  $B$  statistics to estimate confidence intervals or other properties.

### Example R Code for Bootstrapping:

```
1 # Load the boot library
2 library(boot)
3
4 # Create a function for the statistic of interest
5 mean_stat <- function(data, indices) {
6   return(mean(data[indices]))
7 }
8
9 # Sample data
10 data <- c(12.9, 14.2, 13.7, 15.3, 14.8, 13.6, 14.0)
11
12 # Perform bootstrapping
13 bootstrap_results <- boot(data, statistic = mean_stat, R = 1000)
14
15 # Bootstrap confidence interval
16 boot.ci(bootstrap_results, type = "bca")
```

Listing 6: Bootstrap Example in R

### When to Use Bootstrapping:

- When sample sizes are small and standard parametric assumptions are not valid.
- When you need robust estimates of standard errors and confidence intervals.

## 7.5 Choosing Between Advanced Techniques

| Technique             | Data Type               | Purpose                            | R Function                  |
|-----------------------|-------------------------|------------------------------------|-----------------------------|
| Central Limit Theorem | Continuous/Discrete     | Justifies normal approximation     | Manual calculations         |
| Law of Large Numbers  | Continuous/Discrete     | Ensures convergence of sample mean | Manual calculations         |
| Mann-Whitney U Test   | Non-normal, Ranked Data | Compares two independent groups    | <code>wilcox.test()</code>  |
| Kruskal-Wallis Test   | Non-normal, Ranked Data | Compares more than two groups      | <code>kruskal.test()</code> |
| Bootstrap             | Any                     | Estimate confidence intervals      | <code>boot()</code>         |

## 7.6 Key Points

- Advanced techniques are necessary for handling non-standard data and assessing the robustness of models.

- Non-parametric tests are useful when the assumptions of parametric tests are violated.
- Bootstrap methods provide flexible, assumption-free approaches for estimating confidence intervals.

## 8 Useful R Functions and Syntax

R is a powerful language for data analysis and statistical modeling, providing a vast collection of functions for handling, visualising, and modelling data. This section summarises some of the most commonly used R functions to help streamline data analysis.

### 8.1 Data Handling and Manipulation

- `read.csv("file.csv")` — Reads a CSV file into a dataframe.
- `head(df, n)` — Displays the first  $n$  rows of the dataframe.
- `subset(df, condition)` — Extracts rows from a dataframe based on a condition.
- `merge(df1, df2, by = "column")` — Merges two dataframes based on a common column.
- `dplyr` Package Functions:
  - `filter(df, condition)` — Selects rows that meet a condition.
  - `select(df, columns)` — Selects specific columns.
  - `mutate(df, new_var = expression)` — Creates new variables.
  - `group_by(df, column)` — Groups data by a specified column.
  - `summarise(df, new_var = summary_function)` — Aggregates grouped data.

### 8.2 Basic Plotting Functions

- `plot(x, y)` — Generic plotting function for scatter plots, line plots, etc.
- `boxplot(x)` — Creates a box-and-whisker plot for visualising the distribution.
- `hist(x)` — Creates a histogram of a vector  $x$ .
- `barplot(height)` — Plots a bar chart with the specified heights.
- `pairs(df)` — Creates a matrix of scatterplots for pairwise comparison of variables.

### 8.3 Advanced Visualisation with ggplot2

`ggplot2` is a versatile package for creating complex and aesthetically pleasing graphics. It uses a grammar of graphics approach to layer different elements.

- `ggplot(data, aes(x, y))` — Initialises a `ggplot` object with data and aesthetic mappings.
- `geom_point()` — Adds points (scatter plot).
- `geom_boxplot()` — Adds a box plot.
- `geom_histogram(binwidth = value)` — Creates a histogram.
- `facet_wrap(~ factor)` — Creates a series of plots for each level of a factor.
- `theme_minimal()` — Applies a minimal theme to the plot.

### 8.4 Statistical Modelling Functions

- `lm(formula, data)` — Fits a linear regression model.
- `glm(formula, family = "binomial")` — Fits a logistic regression model.
- `anova(model1, model2)` — Compares two nested models using ANOVA.
- `predict(model, newdata)` — Generates predictions based on a fitted model.
- `summary(model)` — Summarises the results of a model fit.

## 8.5 Cheat Sheet Table: Common R Functions

| Category           | Function   | Description                                 |
|--------------------|--|---|
| Data Import        | <code>read.csv("file.csv")</code>                    | Reads a CSV file into R.                    |
| Data Inspection    | <code>head(df)</code>                                | Displays the first few rows of a dataframe. |
| Subsetting         | <code>subset(df, condition)</code>                   | Extracts rows based on a condition.         |
| Merging            | <code>merge(df1, df2, by = "col")</code>             | Merges two dataframes by a common column.   |
| Summarising        | <code>summarise(df, new_var = mean(col))</code>      | Computes summary statistics.                |
| Basic Plots        | <code>plot(x, y)</code>                              | Creates scatter or line plots.              |
| Histograms         | <code>hist(x)</code>                                 | Plots a histogram of the variable $x$ .     |
| Box Plots          | <code>boxplot(x)</code>                              | Creates a box plot to show data spread.     |
| ggplot2            | <code>ggplot(df, aes(x, y))</code>                   | Initialises a ggplot object.                |
| Linear Modelling   | <code>lm(mpg ~ hp, data)</code>                      | Fits a linear regression model.             |
| Logistic Modelling | <code>glm(am ~ hp, data, family = "binomial")</code> | Fits a logistic regression model.           |

## 8.6 R Markdown for Reproducible Reports

R Markdown allows you to create dynamic documents that include code, text, and visualisations. Use the following basic syntax:

- Code chunks: ````{r} ... ````
- Inline R code: ``r variable``.
- Text formatting: **bold**, *italics*, # Heading.
- Export formats: PDF, HTML, Word.

## 8.7 Additional Resources and Tips

- Use the `?function_name` to access documentation for any R function.
- Leverage the `tidyverse` packages for a cohesive data science workflow.
- Use RStudio's auto-complete features to explore function arguments and documentation.

## 9 Quick Overview of Statistical Tests

This section provides a concise summary of the key statistical tests covered in this cheat sheet. It includes information on when to use each test, the hypotheses involved, and the appropriate R functions for implementation.

### 9.1 Choosing the Appropriate Test

The selection of a statistical test depends on the type of data, the number of groups, and the underlying assumptions. Use the following guide to select the appropriate test:

- **Comparing Means:**
  - **One-sample t-test** — Compare sample mean to a known value.
  - **Two-sample t-test** — Compare means between two independent groups.
  - **Paired t-test** — Compare means of paired data (e.g., before/after measurements).
  - **ANOVA** — Compare means across more than two groups.
- **Comparing Proportions:**
  - **Chi-squared Test** — Test for independence in a contingency table.
  - **Fisher's Exact Test** — Use for small sample sizes in categorical data.
  - **Binomial Proportion Test** — Test a specific proportion in a binomial setting.
- **Comparing Distributions (Non-Parametric):**
  - **Mann-Whitney U Test** — Compare ranks between two independent groups.
  - **Kruskal-Wallis Test** — Compare ranks between more than two groups.
  - **Kolmogorov-Smirnov Test** — Compare a sample to a reference distribution or compare two sample distributions.
- **Correlation:**
  - **Pearson Correlation** — For linear relationships with normal data.
  - **Spearman Correlation** — For monotonic relationships or non-normal data.
- **Variance Analysis:**
  - **Chi-squared Test for Variance** — Test for the variance of a population.

### 9.2 Overview of Hypothesis Tests

| Test   | Data Type             | Purpose   | Hypotheses   | R Function                                    |
|--|-----------------------|---|--|---|
| <b>One-Sample t-test</b>                       | Continuous            | Test whether the sample mean differs from a known value                           | $H_0 : \mu = \mu_0 \quad H_1 : \mu \neq \mu_0$     | <code>t.test(x, mu = value)</code>            |
| <b>Two-Sample t-test</b>                       | Continuous            | Compare the means of two independent groups                                       | $H_0 : \mu_1 = \mu_2 \quad H_1 : \mu_1 \neq \mu_2$ | <code>t.test(x ~ group)</code>                |
| <b>Paired t-test</b>                           | Continuous            | Compare the means of the same group at different times                            | $H_0 : \mu_d = 0 \quad H_1 : \mu_d \neq 0$         | <code>t.test(pre, post, paired = TRUE)</code> |
| <b>ANOVA (Analysis of Variance)</b>            | Continuous            | Compare means across multiple groups  | $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$              | <code>aov()</code>                            |
| <b>Chi-squared Test</b>                        | Categorical           | Test the association between categorical variables                                | $H_0 : \text{No association}$                      | <code>chisq.test()</code>                     |
| <b>Fisher's Exact Test</b>                     | Categorical (Small n) | Test association between categorical variables with small sample sizes            | $H_0 : \text{No association}$                      | <code>fisher.test()</code>                    |
| <b>Wilcoxon Rank-Sum Test (Mann-Whitney U)</b> | Non-Normal Continuous | Compare two independent groups using ranks  | $H_0 : \text{Ranks are equal}$                     | <code>wilcox.test()</code>                    |
| <b>Kruskal-Wallis Test</b>                     | Non-Normal Continuous | Compare more than two groups using ranks  | $H_0 : \text{Distribution of ranks is equal}$      | <code>kruskal.test()</code>                   |
| <b>Kolmogorov-Smirnov (K-S) Test</b>           | Continuous/Discrete   | Compare a sample to a reference distribution or compare two distributions         | $H_0 : \text{Distributions are the same}$          | <code>ks.test()</code>                        |
| <b>Binomial Proportion Test</b>                | Binary/Categorical    | Test if the proportion of successes in a population is equal to a specified value | $H_0 : p = p_0 \quad H_1 : p \neq p_0$             | <code>binom.test()</code>                     |
| <b>Chi-squared Test for Variance</b>           | Continuous            | Test for the variance of a population   | $H_0 : \sigma^2 = \sigma_0^2$                      | <code>varTest()</code>                        |

Table 1: Summary of Statistical Tests

### 9.3 Example R Code for Various Tests

Here's a quick overview of how to implement these tests in R:

```
1 # One-sample t-test
2 t.test(data$variable, mu = 20)
3
4 # Two-sample t-test
5 t.test(variable ~ group, data = data)
6
7 # Paired t-test
8 t.test(pre_measure, post_measure, paired = TRUE)
9
10 # ANOVA
11 anova_model <- aov(outcome ~ group, data = data)
12 summary(anova_model)
13
14 # Chi-squared Test
15 chisq.test(table(data$var1, data$var2))
16
17 # Fishers Exact Test
18 fisher.test(table(data$var1, data$var2))
19
20 # Mann-Whitney U Test
21 wilcox.test(group1, group2)
22
23 # Kruskal-Wallis Test
24 kruskal.test(outcome ~ group, data = data)
25
26 # Kolmogorov-Smirnov Test
27 ks.test(sample_data, "pnorm", mean = mean(sample_data), sd = sd(
    sample_data))
28
29 # Binomial Proportion Test
30 binom.test(x = 30, n = 100, p = 0.3)
31
32 # Chi-squared Test for Variance
33 library(EnvStats)
34 varTest(data$variable)
```

Listing 7: Example R Code for Different Statistical Tests