

Quantitative Methods 1 – Textbook Readings Summary

Statistics – a way to get information from data, which can be subdivided into two basic areas: descriptive statistics and inferential statistics

Descriptive statistics – methods of organising, summarising and presenting data in ways that are useful, attractive and informative to the reader

Inferential statistics – methods used to draw conclusions about a population based on information provided by a sample of the population

Descriptive Measures

Population – the entire set of observations or measurements under study

Sample – a subset of observations selected from the entire population

Statistical inference – the process of inferring information about a population

Parameter – a measurement about a population

Statistic – a measurement about a sample

Estimate – an approximate value of a parameter based on a sample

Point estimate – a single value given as an estimate of a parameter (e.g. sample mean)

Confidence interval – the proportion of times that an interval estimate will be correct, if the sampling procedure were repeated a very large number of times

Significance interval – measures how frequently an interval estimate will be correct in the long run

Simple random sample – one in which every member of the population has an equal chance of appearing

Sampled population – the actual population from which a sample has been drawn

Variable – any characteristic that can vary

Data – observations of the variable

Three types of data

1. Numerical (or quantitative) data – observations are real numbers
2. Nominal (or categorical) data – observations are categorical or qualitative
3. Ordinal (or qualitative) data – observations are ranked

Graphical representation of nominal data

Frequency distribution – a table of presenting data and their count in each category or class

Relative frequency distribution – frequency distribution giving the percentage each category or class represents

Bar chart – a graph in which vertical bars represent data in different categories

Pie chart – a circle subdivided into sectors representing data in different categories

Bar charts, pie charts and frequency distributions are used to summarise single sets of nominal (categorical) data. Because of the restrictions applied to this type of data, all that we can show is the frequency and proportion of each category. The type of chart to use in a particular situation depends on the particular information the user wants to emphasise.

Descriptive statistics is concerned with methods of summarising and presenting the essential information contained in a set of data, whether the set be a population or a sample taken from a population.

Random variables and discrete probability distributions

Measures of central location

- **Mean:** the sum of a set of observations divided by the number of them
 - A serious drawback of the **mean** is that it is seriously affected by outliers, which are extreme observations. For this reason, in the existence of outliers, the **median** provides a better measure of central location.
- **Median:** the middle value of a set of observations when they are arranged in order of magnitude
 - When there are a relatively small number of extreme observations, the **median** usually produces a better measure of the centre of the data.
 - The calculation of the mean is **not valid** for ordinal and nominal
 - The **median** is appropriate for ordinal data.
- The **mode**, which is determined by counting the frequency of each observation, is the most appropriate for nominal data. **The mode is the most frequent**

Measures of spread:

- **Deviation** – Difference between an observation and the mean. The sum of deviations is always zero.
- **Standard deviation** is the more useful measure of spread. The standard deviation is the square root of the variance. The standard deviation measure is to be used in conjunction with the mean. The standard deviation is not a single population.
- **Range:** the difference between the maximum and minimum observations.
- The **coefficient of variation** of a set of observations is the standard deviation of the observations divided by the mean. The coefficient of variation is reported as a percentage, which effectively expresses the standard deviation as a percentage of the mean.
- **Outlier** - An observation that is less than $Q1 - 1.5(1QR)$.
- The interquartile range is the difference between the first and third quartiles. It represents the middle 50% of the observations.

Percentile - The value of a variable for which p% of observations are less than that value and (100-p)% are greater than that value.

percentile

$$L_p = (n+1) \frac{p}{100}$$

where L_p is the location of the pth percentile.

Quadrants

- **Q1:** where the first 25% of data lies
- **Q2:** where the first/last 50% of data lies, equivalent to the median
- **Q3:** where the last 25% of data lies

Relationship between two variables

- **Covariance** - A measure of how two variables are linearly related.

- Measured in the **same unit as data but squared** (unit of X times unit of Y = unit of X squared = unit of Y squared)

- **Coefficient of correlation** - A measurement of the **strength and direction of linear relationship between two numerical variables**

- **Unit-less** – takes values from **-1 to +1**
- A value of 0 means no association between the two variables
- A value between 0 and 1 means a positive association
- A value between -1 and 0 means a negative association
- The closer to 1 or -1, the stronger the association
- The closer to 0, the weaker the association
- The correlation is a better measure of the linear association between two variables than covariance because correlation is unit free, and therefore measures the strength of the linear relationship, whereas covariance does not.
- The correlation is based on the covariance

- **Least squares method** - A method of deriving an **estimated line** which best fits the data.

- $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$
- β_0 and β_1 are scaled in terms of the number of observations
- The coefficients β_0 and β_1 are derived using the method of least squares, which minimizes the sum of squared errors (SSE): $\sum_{i=1}^n (y_i - \hat{y}_i)^2$
- When using least squares to estimate the coefficients, the following properties will be minimized:
 - ✓ The sum of residuals = 0
 - ✓ The mean of residuals = 0
 - ✓ The sum of squared residuals is minimized

- **Coefficient of determination** - The proportion of variation in the dependent variable that is **explained** by the independent variable

- Equivalent to the square of the coefficient of correlation

Knowing that an estimator is unbiased means that its expected value equals the parameter it is estimating. The variance of the estimator is to the parameter.

Law of Iteration

- $E[XY] = E[X]E[Y] + \text{COV}(X, Y)$
- $\text{var}(Y | X) = \text{var}(Y) - [\text{var}(Y | X)] + \text{var}[E(Y | X)]$

Expected value: $E[X + Y] = E[X] + E[Y]$

Variance: $\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2\text{COV}(X, Y)$

$\text{var}(X - Y) = \text{var}(X) + \text{var}(Y) - 2\text{COV}(X, Y)$

For independent X and Y then: $\text{var}(X + Y) = \text{var}(X - Y) = \text{var}(X) + \text{var}(Y)$

Hypothesis Testing

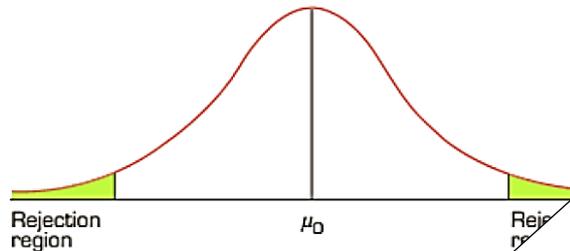
Hypothesis - A proposition or conjecture that the statistician will test by a means called hypothesis testing.

One-tail and two-tail tests

- **Two-tail test:** a test with the rejection region in both tails of the distribution, typically split evenly

$$H_0: \mu = \mu_0$$

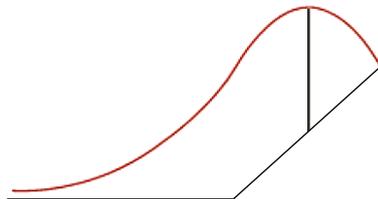
$$H_A: \mu \neq \mu_0$$



- **One-tail test:** a test with the rejection region in only one tail

$$H_0: \mu = \mu_0$$

$$H_A: \mu > \mu_0$$



$$H_0: \mu = \mu_0$$

$$H_A: \mu < \mu_0$$



Right-tail test

Left-tail test

SAMPLE (pg 7)

Power of a test

- The **power of a test** is the probability of rejecting H_0 when it is false = $1 - b = 1 - \text{Pr}(\text{Type 2 error})$
- If you are testing for deviations from the null hypothesis in a single direction, then a one-tail test has more power (less probability of type II error) than a two tail test at the same level of significance. As a result, positive deviations of μ from m have a higher probability of being detected by the upper tail test. So a two tail test is not invalid in this case, but it is not as good as the upper tail test in terms of power

Type 1 and 2 errors

- **Type 1 error:** H_0 is true but is rejected : α
- **Type 2 error:** H_0 is false but is not rejected : b
- As alpha increases, the probability of a type 1 error increases, the probability of a type 2 error decreases, and therefore the power of the test increases.
- As α decreases, the level of confidence for the interval will increase
- Increasing the level of confidence comes by decreasing α , the probability of a type I error. That is, increasing the confidence level is equivalent to saying we want to decrease the probability of rejecting the true value for the parameter under test.