

9. Regression

- used to make predictions on the basis of information that one has obtained
 - Scatterplots can diagnose outlier points
 - x variable \rightarrow horizontal axis \rightarrow predictor (independent variable)
 - y variable \rightarrow vertical axis \rightarrow criterion (dependent variable)
- The regression line is the best fitting line that can be drawn through the points on the scatterplot.
- \rightarrow straight line merely approximates the data it represents \rightarrow it can never pass through every value unless ~~is~~ its correlation is perfect.
 - \rightarrow regression line allows us to make predictions.

Regression equation:

$$\hat{y} = bx + a$$

\hat{y} = predicted y value - value on y axis

b = slope of the line

a = the y intercept point (point where the line intersect y axis)

x = predictor variable (value on the x-axis)

\rightarrow Regression line minimized the squared distances of the observed data points from the line.

Computing the values in the regression equation.

$$b = \frac{N \sum XY - (\sum X)(\sum Y)}{n \sum X^2 - (\sum X)^2}$$

n = number of observation

$\sum XY$ = total of X times Y (remember do each line (score) + then add together)

$\sum X$ = total of X scores/values

$\sum Y$ = total of Y scores/values

$\sum X^2$ = total of X^2 scores (calculate each line and then total)

$(\sum X)^2$ = total of $\sum X$ times by $\sum X$

- If the computed value is

positive = positive slope runs bottom left to top right

negative = negative slope - top left to bottom right.

$$a = \bar{Y} - b\bar{X}$$

\bar{Y} = mean of Y values ($\sum Y/n$)

\bar{X} = mean of X values ($\sum X/n$)

b = slope - (computed above)

∴ can now replace $a + b$ in the equation.

You can now work out the regression line at each point by inserting each X score + calculating.

Drawing the regression line.

- Straight line

2 ways of drawing the line.

Method 1:

- * calculate 2 points
- * mark on scatter plot
- * join 2 points to form a straight line

Method 2:

- * use the Y intercept (a) + any other pair of numbers
$$Y = bx + a$$
- * calculate Y intercept by replacing X with zero in the formula (X = 0 on the X axis)
- * calculate another point by replacing X with an score.

- 'within dataset' predictions are called interpolations
- predictions 'beyond' the range of values in the dataset are called extrapolations
- predictions from regression equations are more accurate when they are made about points that fall within the range of points covered by original dataset.
- Standard error of estimate measures the degree to which the regression lines "fits" the observed data. (Do not need to calculate)

→ The degree of scatter around the straight line is called the correlation

→ Correlation is a measure of the strength and direction of the linear association between 2 variables

→ Casual Inferences should only be drawn from regression analyses with great caution

Standard error of estimate = (know what it is) no need to calculate.

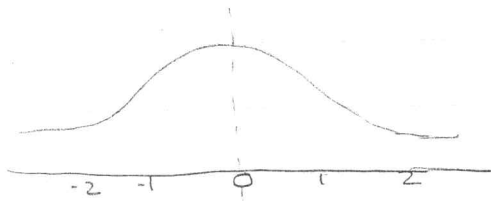
→ if points on the original scatterplot have no clear pattern along with a wide scatter of points, then the line will have very little predictive power + will not be very meaningful.

→ if regression line is a good fit - distances should be quite small - shorter distances from the regression line.

→ cannot simply average the distances → it could balance each other out → by using the standard error of estimate we can take account of ~~the~~ this difficulty (error)

10 Normal distributions

- normal distributions have the following characteristics
- * Bell shape → clearly identifiable bell shape with one distinct peak → unimodal
 - * Symmetry → figure forms an exact mirror image on either side of the center line
 - * area of portion → below the normal distribution equals 100%. Since figure is symmetrical it means that the area on the left occupies 50% and the right the other 50%



Standard normal distribution

- distribution is defined by its shape, mean and variance
- normal distributions allow us to determine where an individual score lies relative to other scores in a set of scores
- infinite number of normal distributions, each with an unique mean, and variance
- to determine proportion of cases falling above or below a score in a distribution we need 3 pieces of information about the distribution
- shape → mean → variance.

→ standard normal distribution is a normal distribution with a mean equal to 0 and a variance equal to 1

→ z scores indicate the number of standard deviation units a score lies above or below the mean.

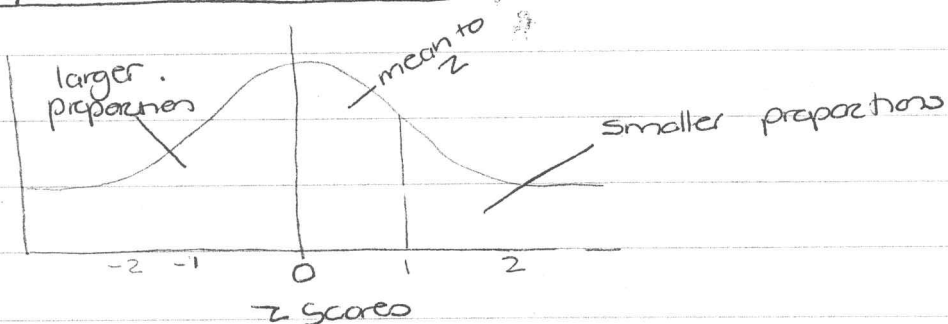
→ standard normal distribution is defined in terms of standard deviation, which are called z scores

→ z-scores range from 3 standard deviations above it and 3 standard deviations below it.

→ The z-table tells us what proportion of the area under the curve of a z distribution lies above or below a particular z score.

* When working with z scores its always best to draw sketch of the distribution. Make the center $z=0$ all negative z scores are on the left and positive on the right

Using the z score tables



→ Although the z-table only gives us proportions lying above positive z scores *

1) because the total area under the curve is 1, we can calculate the proportion below a positive z-score by subtracting the area above the z-score from 1 and

2) because the normal distribution is symmetrical
* the proportion lying above a positive z-score is the same as the area below the negative z-score of the same absolute value.

- The z-table allows us to look up proportions, given z-scores as well as determine z-scores given proportions.

- Standard normal distribution is very convenient since any set of normally distributed scores can be converted to standard scores to enable us to directly compare data from different sets by using

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

- once converted to standard scores, the standard table of z-scores.

→ conversion does not change the distribution of scores only the way they are represented.

⇒ 2 fixed limits commonly used 95% + 99% z-values are.
95% → ±1.96 → 99% → ±2.58

11 Basic concepts of probability

Probability \rightarrow used in the definition of power

\rightarrow no simple definition for probability.

\rightarrow Probability is defined as $p = \frac{a}{n}$

\rightarrow number of successes divided by the total number of events.

\rightarrow events are said to be independent when they do not affect the probability of each other occurring

\rightarrow single event has more than 2 outcomes example!.. pull a card from a deck of cards and you have 52 possible outcomes.

\rightarrow may bet it will be hearts \rightarrow only 13 hearts in the pack \rightarrow the probability of success is thus

$$p = \frac{a}{n}$$

$$= \frac{13}{52}$$

$$= 0.25$$

$$= 1 \text{ in } 4 \text{ or } 25\% \text{ chance / probability.}$$

- Probabilities are based on a past evidence or on the chance of a random event occurring

\rightarrow events must be independent or we will not be able to easily estimate the probability of events taking place.

→ probability law of conjunctions (multiplicative law)

→ probability of 2 independent events jointly occurring is the product of their individual probabilities

$$p(a \text{ and } b) = p(a) \times p(b)$$

Random sampling without replacement undermines independence of events.

∴ when 2 events are not independent from each other we must multiply the probabilities of the outcomes together.

example: tossing a coin twice & getting the same result twice → heads in the first and the second toss will be calculated as follows:

2 possible outcomes: heads
tails

$$\begin{aligned} \therefore p &= \frac{1}{n} \\ &= \frac{1}{2} = 0.5 \\ &= 50\% \text{ chance} \end{aligned}$$

$$\begin{aligned} \therefore (p_a) \times (p_b) \\ &= 0.5 \times 0.5 \\ &= 0.25 \end{aligned}$$

∴ a 25% chance of that happening

The same can be done with drawing a card a second time if card is ~~replaced~~ placed back in card deck.

$$\begin{aligned} \text{it will be } & 0.25 \times 0.25 \\ & = 0.0625 = 6\% \end{aligned}$$

→ if the card is NOT placed back into the deck the deck

is reduced to 51 and the conditions has thus changed \rightarrow
the events are no longer independent to each other and
we have to calculate the new probability $\frac{12}{51} = 0.235$

\Rightarrow now we multiply these probabilities

$$= p(a) \times p(b)$$

$$= \frac{13}{52} \times \frac{12}{51}$$

$$= 0.25 \times 0.235$$

$$= 0.05875$$

$$= 0.059$$

Probability law of disjunction (additive law)

- the probability of either of 2 independent events
occurring is the sum of their individual probabilities

$$p(a \text{ or } b) = p(a) + p(b)$$

Probability of drawing a diamond or a heart

$$= \frac{13}{52} + \frac{13}{52}$$

$$= 0.25 + 0.25$$

$$= 0.5$$

Independent event \rightarrow when the occurrence of an event has no
effect on the probability of an occurrence of
another

Mutually exclusive \rightarrow when the occurrence of one event precludes
the occurrence of another

Exhaustiveness \rightarrow A set of events representing all probable
outcomes

Law of disjunctions (Additive law) \rightarrow events are mutually exclusive.

Law of conjunctions (multiplicative law) \rightarrow when events are independent.

12 Sampling distributions and hypothesis testing.

Sample \rightarrow subset of the population \rightarrow generalize from the sample

It was proven that generalization from a sample is never so different from the result that would be obtained if the entire population average was obtained \therefore it is justified to generalize from sample to a population

Sampling \rightarrow the process of identifying + grouping a subset of an entire population

Random sampling means identifying and collecting the sample in such a way that every member of the population has an equal chance of being in the sample.

\therefore sample is representative of the population

Distributions \rightarrow graphical representation of the ~~entire~~ entire set of data.

Sampling distributions of the mean \rightarrow the frequency distribution of sample means, not individual scores.

\rightarrow one of the main functions of statistics is to make inferences

Drawing inferences involves estimating properties of population from information about samples.

→ many ways to make such estimates → in everyday life people make unscientific inferences → which they treat this hearsay as fact.

→ To draw inferences a representative, random sample is required

→ it is impossible to gain access to the entire population ∴ research relies on drawing of sample cases

Sampling means

when we sample repeatedly from the same population, we expect the means of these samples to be different.

→ Draw several samples from population + get a mean for each sample

→ Now take the mean of all the samples and construct a frequency distribution for these means, we would have a sampling distribution of the mean for the samples

→ Plotting the means of an infinite number of samples of size n , drawn from a population will give us the sampling distribution of the mean.

→ even though the distributions look the same at first glance there are differences

* Scores - different sample values

* Means - means will be different as sample scores are different

* Variances - different samples/means have different variances

Sampling distribution of the mean \rightarrow distribution of an infinite number of sample means of a particular size randomly selected from a population

\therefore impossible to compute the average of an population + therefore we compute the sample mean using the notation \bar{X} \Rightarrow Assume everything that apply to \bar{X} applied to the population μ

\therefore we are ultimately interested in the population we use μ rather than \bar{X}

Hypothesis Testing

hypothesis \rightarrow tentative statement of a relationship between 2 variables

Neuman's definition \rightarrow educated guesses about how the social world works.

hypothesis \rightarrow an idea \rightarrow expectation that there is a correlation between 2 or more variables

Research hypothesis \rightarrow define a statement about the variables
eg: there is a difference between height of women + men in South Africa

How do hypotheses fit into the research process?

\rightarrow have to prove a statement scientifically.

\rightarrow need to follow a process in order to prove or disprove the statement statistically (scientifically)

Sampling distribution of the mean \rightarrow distribution of an infinite number of sample means of a particular size randomly selected from a population

\therefore impossible to compute the average of an population + therefore we compute the sample mean using the notation \bar{X} \Rightarrow Assume everything that apply to \bar{X} applies to the population μ

\therefore we are ultimately interested in the population we use μ rather than \bar{X}

Hypothesis Testing

hypothesis \rightarrow tentative statement of a relationship between 2 variables

Neuman's definition \rightarrow educated guesses about how the social world works.

hypothesis \rightarrow an idea \rightarrow expectation that there is a correlation between 2 or more variables

Research hypothesis \rightarrow define a statement about the variables
eg: there is a difference between height of women + men in South Africa

How do hypotheses fit into the research process?

\rightarrow have to prove ~~it~~ a statement scientifically.

\rightarrow need to follow a process in order to prove or disprove the statement statistically (scientifically)

9 steps in statistical / hypothesis testing

- 1) * formulate a null hypothesis \rightarrow a statement that maintains that there is no difference between the groups or conditions. It is represented by the symbol H_0
- \rightarrow negative statement about a positive difference between 2 or more variables

eg: There is no difference between the height of men + the height of women

OR Men and women are equally tall

Formulated as: $H_0 = \mu_m = \mu_w$

OR $\mu_m - \mu_w = 0$

- 2) Formulate the alternative hypothesis

\rightarrow affirmative statement. \therefore there is a difference between variables.

* non-directional \rightarrow 2 variables are not equal. \Rightarrow there is a difference between the height of women + men

* directional \rightarrow one variable is smaller than the other \Rightarrow men are shorter than women

* directional \rightarrow one variable is larger than the other \Rightarrow men are taller than women

example: There is a difference between the job performance between women + men

$H_1: \mu_m \neq \mu_w$

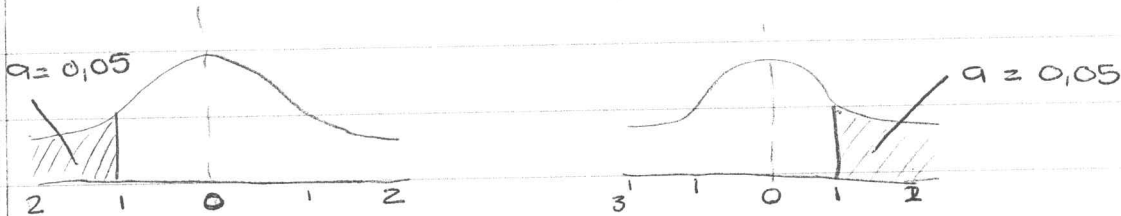
$H_1: \mu_m - \mu_w \neq 0$

3) Determine whether the test is one-tailed or two-tailed.?

→ one tailed → directional → one tailed < or >

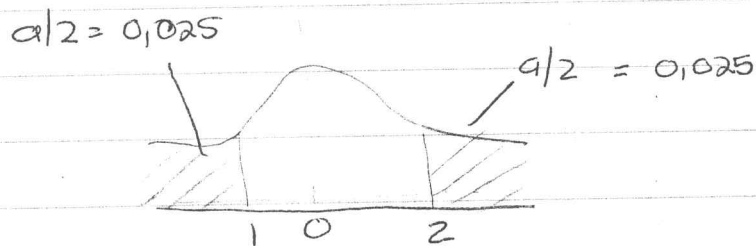
→ two tailed → non directional → not sure then =

Directional - one tailed < or >



Rejects extreme outcomes at only one end of the distribution

Non-Directional - two tailed



4) Determine the level of significance.

work with 5% or a 1% level of significance

Type I error ⇒ rejecting a null hypothesis when it's true

Type II error ⇒ not rejecting a null hypothesis when it's false

probability with which we reject are willing to reject the null hypothesis when it's correct

5) Compute the test statistic.

→ use a few test statistics

T test

F test (chi-square test)

when we want to determine whether there is a significant difference between sample means

→ Test statistics in which the mean plays a role can be classified in 2 categories: those relating to only 2 samples (t-tests) and those relating to more than 2 samples (F test)

2 samples → t test for related samples

t test for unrelated samples

3 or more samples → F test (one way analysis of variance)

1 or more variables with frequencies - chi-square test

Test statistics F + Chi-square refer to the answers that you obtain when testing a hypothesis empirically (statistically) and to where the calculating of the test-statistics fit into the 9 steps of hypothesis testing.

6) Determine the degrees of freedom

→ need to determine the degrees of freedom as to interpret your computations properly.

7) Determine the critical value.

take into account

* whether it is a one tailed or 2 tailed test.

* level of significance

* degrees of freedom

to determine critical value.

critical value is an indication of where rejection region begins + indicates that all scores above or below this point fall into the 5% or 1% level of significance.

8) Reject or do not reject the null hypothesis

* if test statistic $>$ critical value reject H_0

* if test statistic $<$ critical value do not reject H_0

~~Always~~ Always work with absolute values \rightarrow ignore the negative sign + read as positive.

9) Interpret the finding

- Interpret the rejection or non rejection of the null hypothesis in relation to the original research problem.

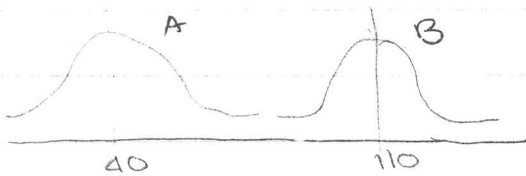
NB Make sure you can write down the 9 steps

13 hypothesis tests applied to means t -tests

t -test is used to compare 2 (estimated) population means

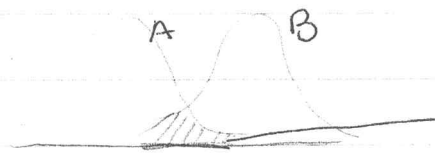
→ compare normal distributions

How do you know if the distributions is from the same population.



A - lower scores
B - higher scores

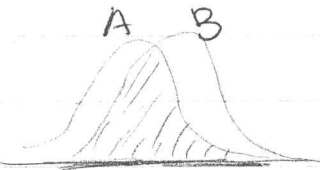
∴ not overlapping - scores from A unlikely to appear in B



Distributions overlap.

Scores can belong to A or B

Separate but similar as they do overlap a number of scores.



- overlap remarkably
- appear the same

T -test allows us to determine the degree to which 2 distributions overlap.

The degree to which 2 distributions overlap is determined by the difference between the means and variance of each sample.

$$t = \frac{\text{difference between the means}}{\text{Standard error}}$$

$$z = \frac{\text{difference between score + population mean}}{\text{Standard deviation of the population}} = \frac{x - \mu}{\sigma}$$

Central difference between t-test and the z-score transformation is that the population parameters ($\mu + \sigma$) do Not have to be known to perform a t-test.

This is important because in the large majority of cases the population parameters are unknown.

⇒ unlike a z-test no population parameters are needed to perform a t-test

Standard error is the standard deviation of the sampling distribution of the mean (the distribution created by taking repeated sample means from a population)

$$S_x = \frac{S}{\sqrt{n}}$$

To analyse data with a t-test, the data need to comply with the assumptions of normality, homogeneity of variance and independence.

purpose of hypothesis testing with 2 independent means is to help you decide whether an observed difference between 2 sample means is accidental or whether it represents a real difference between populations

- Independent samples

* Sampling distribution of difference between means is the distribution of an infinite number of a particular size randomly selected from a population

* because we sample each population independently, the sample means will also be independent

* mean of the sampling distribution will be $\mu_1 - \mu_2$

Formula used for t -value of 2 independent Groups.

* requires the use of the variance of both groups

if you do not have the variance, you must compute it

using:

$$s^2_x = \frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n-1}$$

Do the same for y (second variable)

Now compute the t -value

Steps 1-9 (Chapter 12)

- Formulate the null hypothesis
- Formulate the alternative hypothesis
- Determine if the test is 1 or 2 tailed
- Determine level of significance
- Compute test statistic

$$= t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{S_1^2}{N_1} + \frac{S_2^2}{N_2}}}$$

- Determine degrees of freedom

$$df = N_1 + N_2 - 2$$

- Determine critical value (table A1.2)
- Determine the null hypothesis should be rejected/not rejected
- Interpret the finding

Pooled variances \rightarrow when we deal with samples of unequal sizes

- Use the pooled variance as your standard error estimate unless the sample variances are not homogeneous

$$S_p^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

n_1 = sample size of 1st sample

n_2 = sample size of 2nd sample

S_1^2 = variance of first sample

S_2^2 = variance of second sample.

T-test Independent Samples

$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_{\bar{X}_1 - \bar{X}_2}} \quad \text{or} \quad t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{S_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

\bar{X}_1 = mean of 1st Sample

\bar{X}_2 = mean of 2nd sample

$S_{\bar{X}_1 - \bar{X}_2}$ = estimate of the standard error.

use separate variance estimates if the variances of your 2 samples are highly different (one variance is at least 4 times as big as the other)

$$S_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

S_1^2 = variance 1st Sample

S_2^2 = variance 2nd sample

n_1 = sample size of 1st sample

n_2 = sample size of 2nd sample

heterogeneity variance \rightarrow when samples are drawn from the populations have different variances

\rightarrow we compensate for heterogeneity by computing the t-test separately for the different variances (not pooling them)
smaller n is then used to read the degrees of freedom.

- 13 Non-parametric equivalent of the t-test
- statistical tests that do not rely on parameter estimation and/or distributional assumptions
 - ↳ Mann-Whitney test.
 - ↳ tests whether 2 independent samples have the same median

Wilcoxon matched pairs test → tests whether 2 related samples have the same median

Wilcoxon → signed ranks test
→ rank order differences (ignore sign of difference)

Related Samples

- ↳ samples that are dependent.
- ↳ 2 sets of data on the participants (same participant)
- ↳ we expect data to correlate.

∴ 2 groups of people / participants are matched according to a variable in 1 group. → Group 2 assigned a member that corresponds with the member of group 1 in respect of the particular variable.

∴ Particular pairs are selected in respect of the relevant variables → age, height, eye color can be used to match the persons / variables

2 ways of obtaining related sample scores

→ by matching

→ when 1 participant contributes 2 scores

- * Compute the mean \bar{D}
- * standard deviation ($s.D$) of different scores

t-test for related groups

$$t = \frac{\bar{D} - 0}{\frac{SD}{\sqrt{N}}}$$