



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence



ECON2271

Business Econometrics

Or: Practical Econometrics for Beginners

Week 1 Workbook Notes*

* Workbook notes are partially complete lecture notes intended to trick your brain into action and stop your eyes from glazing over. Take these notes to the lecture and complete the material in class.

Week 1: Introduction to key concepts

Key Unit Information



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Formal teacher-led activities:**

1. Two-hour (90-minute) workshop (lecture), Wednesday 10am
 - Goes through all key concepts with examples and exercises
 - Intended to be an interactive workshop rather than a traditional lecture
 - Students must download and print incomplete slides prior to class, to be completed in class
2. One hour (45-minute) lab in the Trading Room
 - Starts in week 2, and run most weeks
 - Students learn how to use Stata to work with data sets
 - Students given exercises to perform and interpret
 - Exercises are assessed!

Week 1: Introduction to key concepts

Key Unit Information



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Topics covered:**

1. Topic 1: Key concepts: Data, distributions, and comparing means
2. Topic 2: Univariate regression with continuous and unlimited depended variables
3. Topic 3: Multivariate regression with continuous and unlimited depended variables
4. Topic 4: Binary, limited and discrete dependent variable models
5. Topic 5: Extending data formats beyond cross-sections: Time series and Panel data
6. Topic 6: Research design (time permitting)

CRCOS Provider Code: 00126G

Week 1: Introduction to key concepts

Key Unit Information



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Assessments:**

1. **Lab activities:** (about 8) **20%**
 - Focus on ability to interpret results obtained using Stata
 - Assessed component is completed in pairs
 - Students must circulate – not allowed to have the same partner
2. **Tests (2):** **30%**
 - Test 1: 15% Topics 1, 2, Friday 24 August at 4pm
 - Test 2: 15% Topic 3, Friday 21 September at 4pm
3. **Final exam:** **50%**
 - Covers all material

CRCOS Provider Code: 00126G

Week 1: Introduction to key concepts

Key Unit Information



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- Unit structure and proposed schedule:

Week	Date	Workshop	Lab	Test
1	30-Jul	Topic 1		
2	6-Aug	Topic 2 (i)	Lab 1	
3	13-Aug	Topic 2 (ii)	Lab 2	
4	20-Aug	Recap		Test 1
5	27-Aug	Topic 3 (i)	Lab 3	
6	3-Sep	Topic 3 (ii)	Lab 4	
7	10-Sep	Topic 3 (iii)	Lab 5	
8	17-Sep	Recap		Test 2
MSB	24-Sep			
9	1-Oct	Topic 4	Lab 6	
10	8-Oct	Topic 5 (i)	Lab 7	
11	15-Oct	Topic 5 (ii)	Lab 8	
12	22-Oct	Topic 6, recap		

CRCOS Provider Code: 00126G

Week 1: Introduction to key concepts

Key Unit Information



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- Key resources:

- Unit LMS page: One-stop-shop for all your resources
- Lecture notes (incomplete): Released weekly (as required)
- Exercises and lab instructions: Released weekly (as required)
- Discussion board: For Q&A, comments and feedback
- Recommended text: Gujarati's Basic Econometrics or equivalent

CRCOS Provider Code: 00126G

Week 1: Introduction to key concepts

Agenda and learning outcomes



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Data, distributions and comparing means**
 - Data scales and quality
 - Students able to identify data scales and quality
 - Cross-sectional, time-series and panel
 - Students able to identify format and implications for analysis
 - Distributions, sampling and distributional characteristics
 - Students able to understand the concept of a probability and frequency distribution and sampling, key characteristics of the normal distribution, and key measures of central tendency and dispersion.
 - Testing means: hypotheses, tests and confidence intervals
 - Students able to identify and test hypotheses comparing means “manually”, identifying and interpreting a 95% confidence interval for a mean.

CRCOS Provider Code: 00126G

1. Types of Data:

Informational quality



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- Data are categorized by the length and distinctness of the measurement scale:
 1. **Binary “scale” (categories):**
 -
 - E.g. Are you currently in paid work? (yes/no)
 2. **Discrete scale:**
 -
 - E.g. How many children do you have? (only whole numbers possible)
 3. **Continuous scale:**
 -
 - E.g. Length of a piece of string
 4. **Unlimited scale**
 -
 - E.g. net worth (can extend either end of zero, with no limit)
 5. **Limited scale:**
 -
 - E.g. Length of a piece of string and number of children are bounded at one end (zero); a 0-100 index is bounded at both ends.

CRCOS Provider Code: 00126G

1. Types of Data:

Informational quality

- Data are categorized by their “informational quality”:

1. Nominal scale:

-
- E.g. Transportation type: can be car, bus, train, walk, cycle...
- Does not contain any information beyond distinctness, such as order.

2. Ordinal scale:

-
- E.g. Credit ratings: AAA, AA, A, B...
- Contains information on distinctness and order: We can say that $AAA > AA > A$ etc.
- But does not contain information beyond order, e.g. how much better is AAA compared to AA, and AA compared to A, etc.
- So makes no sense to add or subtract these data.

1. Types of Data:

Informational quality

3. Interval scale:

-
- E.g. temperature in degrees C, calendar date
- Contains information on distinctness, order and also distance between values: We can say that the temperature today is 20 degrees, which is 3 degrees warmer than today; and that Jack and Jill were born 3 days apart.
- But lacks a non-arbitrary (definitive) zero-point. Celsius and Fahrenheit temperatures have arbitrary zero-points (though you might argue that Kelvin degrees do not); and the same is true for calendar.
- This means ratios have no meaning

4. Ratio scale:

-
- E.g. dollars, measurements of length, volume, weight, displacement etc.
- Contains information on distinctness, order, distance between values AND distance from a non-arbitrary (definitive) zero-point.
- These data have complete informational quality, such that ratios are meaningful (as well as sums and differences).

1. Types of Data:

Data formats



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

1. Cross-sectional data:

-
- E.g. the height of 3 year-old girls in a community, the GDP of a set of countries, the incomes of households in a suburb, the profits of firms in an industry, or the marks for a class of ECON1101 students.
- These data have one dimension, which is the unit of interest: which can be 3-year-old girls, countries, households, firms or students in ECON1101.
- We can observe variation across units.

2. Repeated cross-sectional data:

-
- E.g. the height of 3-year-old girls in a community in 1980, 1990, 2000, and 2010; the marks for ECON1101 students across semesters.
- We can observe variation across units, and compare the characteristics of these data sets at different points in time. E.g.: Are 3-year old girls getting taller over time? Are my ECON1101 students doing better this semester, compared to last semester?

CRCOS Provider Code: 00126G

1. Types of Data:

Data formats



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

3. Time-series data:

-
- E.g. the price of Woolworth shares, Australian GDP, the RBA cash-rate, the value of Australian exports to China, my consumption of coffee...
- These data have one dimension, which is time.
- We observe variation across time, for a given unit.
- Frequencies differ: share prices change minutely (perhaps even more frequently?), GDP is calculated quarterly (I think), the cash rate is set monthly, my consumption of coffee changes daily...

4. Panel data:

-
- E.g. Annual financial year incomes of a set of households.
- These data have two dimensions: units of observation and time
- These data sets can be very powerful and allow us to discover things we cannot observe with other data formats. More on this later in the unit...

CRCOS Provider Code: 00126G

2. Data distributions

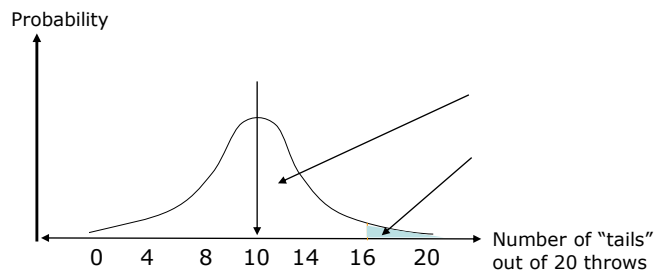
Probability (frequency) distributions



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

Probability (frequency) distributions:

- Randomness is everywhere: I know that there is a 50% chance of getting “tails” or “heads”, respectively, every time I toss a coin, but if I toss the coin 20 times, I can’t be 100% certain that I’ll get “tails” 10 times and “head” 10 times. However, this is the most likely outcome. The second most likely outcome is 9 heads and 11 tails, or 11 heads and 9 tails; followed by 8+12 and 12+8, etc. It is even possible that I’ll get 20 tails or 20 heads – though it’s very unlikely. In fact, there is a probability distribution that tells me how likely each different outcome is:



CRICOS Provider Code: 00126G

14

2. Data distributions

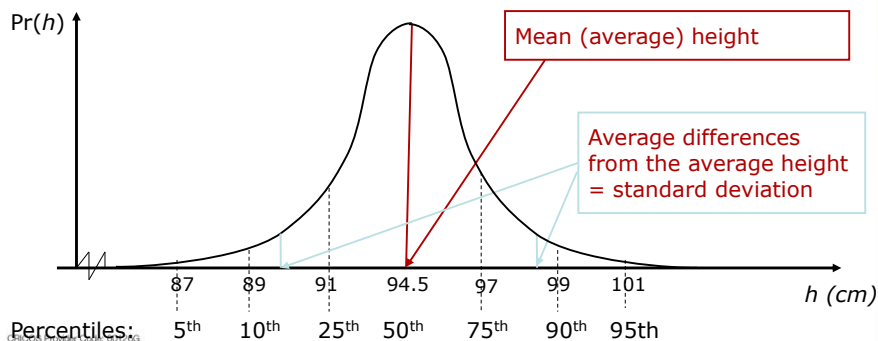
Probability (frequency) distributions



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

The normal distribution:

- A probability distribution maps the probabilities of observing a particular random event (which ideally is measured on a continuous scale between $-\infty$ and ∞ , though this is rarely the case in reality).
- Example of a normal distribution:** The height of 3 year-old girls:

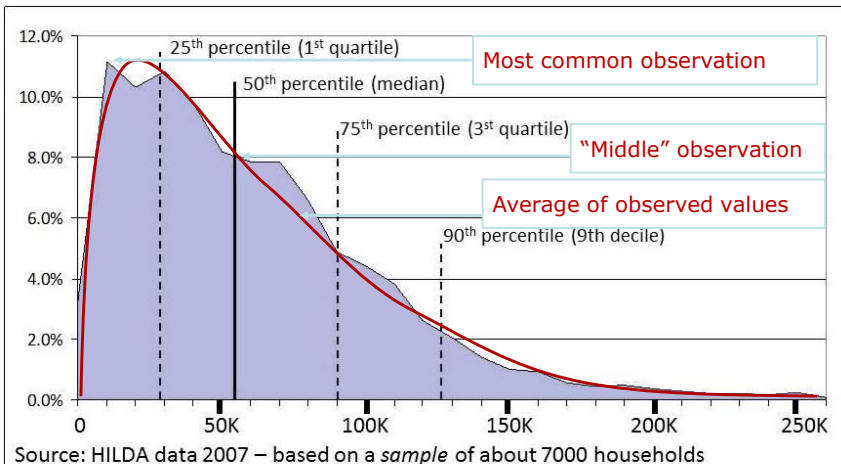


2. Data distributions

Probability (frequency) distributions

- **A non-normal distribution:**

- E.g. Household incomes in Australia (disposable after-tax):



2. Data distributions

Sampling

- **Samples versus populations:**

- A sample is a selection drawn from a population
 - E.g. if I look at a selection of 200 students of the 1000 enrolled in ECON1101, in order to tell me something about people enrolled in this unit, then the 200 students represent a sample, and the 1000 students in total represents the population.
 - E.g. if I look at the students enrolled in ECON1101 in order to tell me something about first year students, then the ECON1101 students is my sample, and all first year students is the population.
 - E.g. I find that only 2 out of 15 mothers in a mothers' group use biodegradable nappies (for their children!!). I conclude that 13% of households with babies use biodegradable nappies. I then use the mothers' group as a sample and use their behaviour to infer something about the entire population of households with babies.
- Is that reasonable?
- How large should a sample be in order for me to be able to make such inferences?

2. Data distributions

Sampling



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Samples versus populations:**

- When I find that the pass rate in ECON1101 is 77% in a particular semester, this is not an inference. It's the truth! 77 out of 100 students pass in this semester. No ifs or buts.
 - I am looking at the population and I know "everything" about it!
 - There is no "inference" involved
- When I look at 10 representative undergraduate units which have an average pass rate of 76%, then I can use this to make an inference about the general pass rate of undergraduate units. This doesn't mean that all units have this pass rate, and I am just hoping that this selection is representative and not too different from other units. If my selection is representative, my average pass rate is a good indication of what's going on in undergraduate units generally.
 - I am looking at a sample, which I know everything about, to make inferences (guesses) about a population which I don't know everything about.

CRCOS Provider Code: 00126G

2. Data distributions

Sampling



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Samples versus populations:**

- **Sample size:**
 - If you sample 8 random students about their study habits, how powerful will your results be?
 - What if you sample 20? 30? 600? 80,000?
- In general, you need at least 30 observations to get results that are of any use, and the greater the sample size, relative to the population size, the more accurate will be the statistical inferences (guesses)
- **Representative samples:**
 - If you are doing market research for a pub where about 70% of clients are male, and you want to collect a sample of pub goers to study their preferences in terms of beer, how useful will a sample with 50% males be to you?
- The more representative a sample is of the general population, the more accurate will be the statistical inferences (guesses)

CRCOS Provider Code: 00126G

2. Data distributions

Sampling



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Samples versus populations:**
 - If we are studying a population distribution, there is no guesswork involved: we know all the data.
 - If we are studying a sample distribution in order to understand the population, we use using statistical inference to provide the best possible guess about that population.
 - Hence, the sample mean and distribution hopefully reflects the population, if the sample is sufficiently large and representative – a true random selection of the population.
 - This means that any number of samples can be drawn from the population, providing a set of sample means. Our chosen sample mean is just one within a whole distribution of sample means
 - However, as the sample size is larger and larger, approaching the size of the population itself, the sample mean will naturally approach the population mean.
 -

2. Data distributions

Distributional characteristics



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Measures of central tendency:**
 - **The mean:**
 - The average observed value
 - Formula:
$$\bar{x} = \frac{\sum x}{n}$$
 - **The mode:**
 - The most frequently observed observation
 - = the value you are most likely to observe
 - **The median:**
 - The “middle” observation
 - = the one which separates the top and bottom 50% of observed values
 - When sorting all observations from lowest to highest (or vice versa), the median is the one in the middle.

2. Data distributions

Distributional characteristics



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Measures of central tendency:**
 - **Example:** Data on wealth across households result in the following statistics:
 - Wealth per household (=assets – debt):
 - Mean: \$604,450
 - Median: \$342,000
 - Mode: nil
 - $n = 12,000$
 - Min = -\$916,783
 - Max = \$12,789,073
 - What can you surmise from this?

2. Data distributions

Distributional characteristics



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

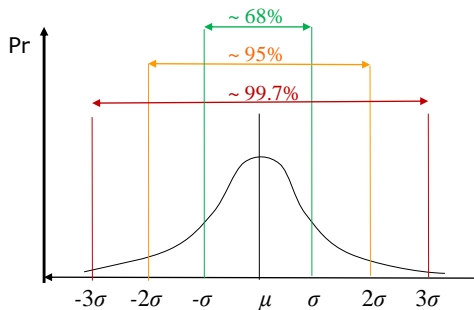
- **Measures of dispersion (variation):**
 - As shown in the previous examples, measures of central tendency are much more informative when we also have some measures of dispersion
 - While measures of central tendency tells us what the average, most likely or middle observation is, measures of dispersion gives information about the variation in the data
 - We have several measures of dispersion:
 - **Minimum and maximum:** The smallest and largest observed values. Simple and very informative.
 - **Range:** The difference between the largest and smallest observed values
 - **Standard deviation:** The average distance away from the mean in the data = the square root of the variance (Variance = σ^2 , s^2):

$$\text{Population : } \sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} \quad \text{Sample : } s = \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$$

2. Data distributions

Distributional characteristics

- **Standard deviations and the probability distribution:**
 - A normal distribution has specific properties:
 - 99.7% of observations will lie within 3 standard deviations (σ 's) either side of the population mean (μ)
 - 95% of observations will lie within 2 standard deviations (σ 's) either side of the population mean (μ)
 - 68% of observations will lie within 1 standard deviation (σ 's) either side of the population mean (μ)



CRICOS Provider Code: 00126G

3. Hypothesis testing

Tests of means

- **The scientific method: Hypothesis testing**
 - We can use samples to estimate tendencies and probabilities (e.g. how likely are women to develop skin cancer over their life time?) as well as associations (e.g. what is the relationship between income and life satisfaction?), and test hypotheses (e.g. are women more or less likely to develop skin-cancer, compared to men?).
 - We can then articulate a null hypothesis, H_0 , which is basically “nothing interesting is going on” (e.g. women and men are not different), and test this against the alternative hypothesis, H_a (e.g. women and men are different).
- **Standard Deviation and Standard Error**
 - We use sample means and standard deviations to evaluate the probability of being wrong in concluding that the population mean is equal to a given value, or different from a given value; and that two sample distributions are distinct or equal.
 - When comparing sample means we use standard errors, $SE = \frac{s}{\sqrt{n}}$
 - While the standard deviation measures the amount of variability or dispersion from the mean, the SE measures how far the sample mean of the data is likely to be from the true population mean. SE is always smaller than s .

CRICOS Provider Code: 00126G

3. Hypothesis testing

Tests of means



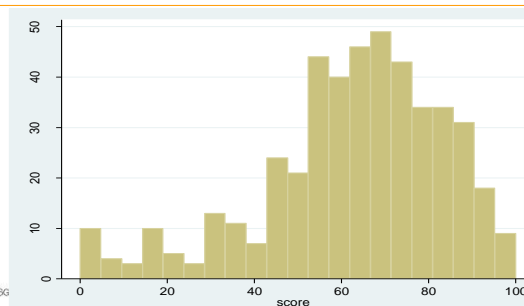
THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Example 1: Final marks for ECON1101**

- The final marks for ECON1101 may be treated as a population. In that case, no guesswork is required. These are the final marks for the unit in semester 2, 2017 (here using Stata, though you can also find the same statistics easily using Excel):

```
. summarize score
```

Variable	Obs	Mean	Std. Dev.	Min	Max
score	459	62.56645	21.2223	0	100



CRCOS Provider Code: 00126G

3. Hypothesis testing

Tests of means



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Example 1: Scaling of final marks for ECON1101**

- The final marks for ECON1101 may also be treated as a sample of final marks for large first-year units. We can't expect all units to have exactly the same distributional characteristics. If we collect the marks for all large undergraduate units, we get a "distribution of distributions".
- Let's say I am told the mean score for this unit should be about 65/100. If the unit mean score is statistically different from 65, then I need to apply scaling. The question is: is my sample sufficiently different? Does a mean of 65 fall inside of the, say, 95 percent confidence interval for my mean? I.e., is it statistically different from 65 at the 95% level of confidence?
- The null hypothesis to be tested is:
 - H_0 : There is no statistical difference between my mean and the value of 65. This difference is just caused by randomness.
- The alternative hypothesis is:
 - H_a : My mean is statistically different from 65. This difference reflects a fundamental distinction between my distribution and the given population mean, rather than randomness.

CRCOS Provider Code: 00126G

3. Hypothesis testing

Tests of means



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- Testing the sample mean against a given value

- The appropriate test statistic is the t-statistic:

$$t = \frac{\bar{x} - \mu}{SE(\bar{x})} = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

- Here:

- I need to choose how confident I want to be in rejecting the null and accepting the alternative hypothesis. The default appears to be 95%.

```
. ttest score==65
```

One-sample t test

Variable	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
score	459	62.56645	.990572	21.2223	60.61982	64.51308

mean = mean(score)

t = -2.4567

Ho: mean = 65

degrees of freedom = 458

Ha: mean < 65

Ha: mean != 65

Ha: mean > 65

Pr(T < t) = 0.0072

Pr(|T| > |t|) = 0.0144

Pr(T > t) = 0.9928

3. Hypothesis testing

Tests of means



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- Testing the sample mean against a given value

–

- What if my sample size was smaller? I select the first 100 observations and pretend this is my sample instead: (I call this variable "score2")

```
. ttest score2==65
```

One-sample t test

Variable	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
score2	100	61.26	2.279758	22.79758	56.73647	65.78353

mean = mean(score2)

t = -1.6405

Ho: mean = 65

degrees of freedom = 99

Ha: mean < 65

Ha: mean != 65

Ha: mean > 65

Pr(T < t) = 0.0520

Pr(|T| > |t|) = 0.1041

Pr(T > t) = 0.9480

3. Hypothesis testing

Tests of means

- **Example 2:** Comparing ECON1111 marks across two different semesters
 - I ran ECON1111 in semester 2 2017, and semester 1 2018. I found the latter group did much better than the former group, with an average score of about 10 marks higher!! But was this difference in means statistically significant, or might it be put down to randomness?

```
. summarize s22017 s12018
```

Variable	Obs	Mean	Std. Dev.	Min	Max
s22017	61	57.2623	24.0935	3	92
s12018	179	67.51397	17.69033	0	95

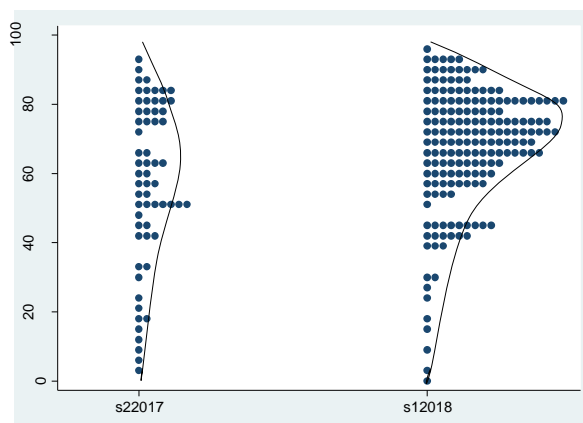
- The null hypothesis to be tested is:
 - H_0 : There is no statistical difference between these two means. This difference is just caused by randomness. I.e., there is substantial overlap between these distributions.
- The alternative hypothesis is:
 - H_a : These two means are statistically different. This difference is random.

CRCOS Provider Code: 00126G

3. Hypothesis testing

Tests of means

- **Example 2:** Comparing ECON1111 marks across two different semesters



CRCOS Provider Code: 00126G

3. Hypothesis testing

Tests of means



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- Testing for differences in means across samples

- Null hypothesis:

- t statistic:
$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} =$$

- Critical values for t-statistic for $\bar{x} = 0$ or $\bar{x}_1 - \bar{x}_2 = 0$ [$\bar{x} > \text{or} < 0$ or $\bar{x}_1 - \bar{x}_2 > \text{or} < 0$]; $n = 60$:
 - 90% level of confidence: $t = 1.671$ [1.296]
 - 95% level of confidence: $t = 2.000$ [1.671]
 - 99% level of confidence: $t = 2.660$ [2.390]

3. Hypothesis testing

Tests of means



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- Testing for differences in means across samples

- You can perform this test easily, using Stata:

```
. ttest s12018 == s22017, unpaired unequal
```

Two-sample t test with unequal variances

Variable	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
s12018	179	67.51397	1.322237	17.69033	64.90469	70.12324
s22017	61	57.2623	3.084857	24.0935	51.09166	63.43293
combined	240	64.90833	1.289075	19.97027	62.36893	67.44773
diff		10.25167	3.356285		3.576304	16.92704

```

diff = mean(s12018) - mean(s22017)          t = 3.0545
Ho: diff = 0                                Satterthwaite's degrees of freedom = 83.1254

Ha: diff < 0                                Ha: diff != 0                                Ha: diff > 0
Pr(T < t) = 0.9985                          Pr(|T| > |t|) = 0.0030                          Pr(T > t) = 0.0015

```


3. Hypothesis testing

Tests of means

- **Confidence Intervals:**

- The sample mean provides an estimate for what the true population mean (which is not known) is. The SE of the sample mean provides an “error margin” around this estimate, which determines how precise this estimate is. A larger SD (or s) means a larger SE, while a larger n means a smaller SE.
- The width of the confidence interval (CI) for the mean is determined by the SE of the mean, and the level of confidence. The default is 95%, usually, which is the mean \pm 2 SEs.
- In this case, I am interested in whether the difference in the two means may be put down to randomness, or whether it is likely that these samples are really drawn from two distinct populations.

Variable	Obs	Mean	Std. Err.	Std. Dev.	[95% Conf. Interval]	
s12018	179	67.51397	1.322237	17.69033	64.90469	70.12324
s22017	61	57.2623	3.084857	24.0935	51.09166	63.43293

- **Interpretation of p-values**

- Please read: <https://www.students4bestevidence.net/p-value-in-plain-english-2/>

CRICOS Provider Code: 00126G

Test your understanding

Week 1 material

- **Exercise 1: Testing of means**

- A large sample of Australian individuals contain information on life satisfaction, which is a numeric score between 0 (=totally dissatisfied) and 10 (=totally satisfied); and a range of other characteristics. A comparison of life satisfaction scores across individuals who were born in an English-speaking country and those from non-English-speaking countries reveal the following statistics:

	n	Mean	S.D.
English-speaking COB	15,395	7.939	1.438
Non-English-speaking COB	2,102	7.693	1.568

Q: Do people born in non-English-speaking countries exhibit lower life satisfaction compared to people born in English-speaking countries?

- State the hypothesis to be tested.
- Identify the appropriate test statistic, and calculate the t-stat.
- Determine whether you can reject the null hypothesis at the 95% level of confidence.

CRICOS Provider Code: 00126G

Test your understanding

Week 1 material



THE UNIVERSITY OF
WESTERN AUSTRALIA
Achieve International Excellence

- **Exercise 1: Testing of means**
 - i. State the hypothesis to be tested.
 - ii. Identify the appropriate test statistic, and calculate the t-stat.
 - iii. Determine whether you can reject the null hypothesis at the 95% level of confidence.