

Edge Hill University

# Advanced Data Science

Topic 11b – Part 5

---

---

---

---

---

---

---

---

Edge Hill University

## 1. What We'll Cover

This topic will introduce...

- What is data science.
- Key concepts – the scientific method.
- Useful terminology.
- Important tools - Statistics.
- Data collection & Experiment Design.
- Probability basics.
- Data distributions.
- Hypothesis testing.

} Part 5

The aim: to help you understand what it means to be a data scientist and to get you familiar with data science tools.

---

---

---

---

---

---

---

---

Edge Hill University

## 2. Hypothesis Testing

- Hypothesis testing is a statistical approach used to find the optimal answer to the questions we pose about the world around us.
- It uses available knowledge captured in data, to reach conclusions regarding hypotheses in a rigorous way.
- The method is useful when undertaking experimental studies.
- Suppose we are tasked with determining if a medicine works.
- We form hypotheses and split a sample population into control and experimental groups.
- We can use hypothesis testing to determine which of the hypotheses holds over the groups.
- That is, which has the most evidence in it's favour.
- Here we are introducing the foundations of statistical inference central to data science and machine learning.

Null Hypothesis

$$H_0 = \text{No effect}$$

Alternative Hypothesis

$$H_a \text{ or } H_1 = \text{Effect}$$

```

graph TD
    A[Population Sample] --> B[Control Group]
    A --> C[Experimental Group]
  
```

---

---

---

---

---

---

---

---

### 3. How this fits in?

- So far you've come across concepts from lots of different areas.
- You've learned about probability theory, the different types of data distribution (unimodal, bi-modal, multi-modal), the law of large numbers, and how to compute summary statistics over samples of data, and entire populations.
- We covered this material to help prepare you for the concepts I'll very shortly introduce related to hypothesis testing. I'm sure you're relieved that none of this time was wasted!
- So with that in mind, lets return to thinking about a distribution I've mentioned a few times during this course – the normal distribution.

---

---

---

---

---

---

---

---

### 4. Back to Normal

- The normal curve is always a symmetric, unimodal, bell-shaped curve.
- The shape of the curve is determined by two parameters.
  - The mean,  $\mu$ .
  - The standard deviation,  $\sigma$ .
- We can describe any normal curve via a pair e.g.  $(\mu = 0, \sigma = 1)$ .

Standard Normal Distribution

---

---

---

---

---

---

---

---

### 5. Back to Normal

---

---

---

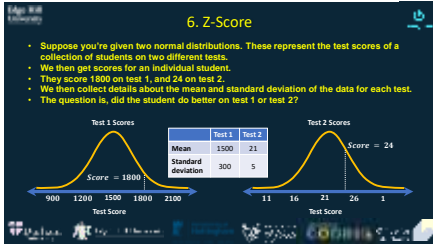
---

---

---

---

---




---

---

---

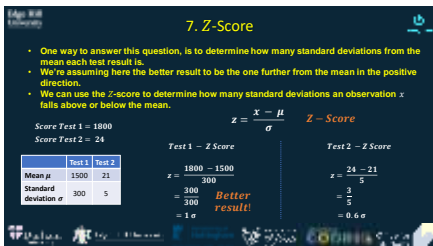
---

---

---

---

---




---

---

---

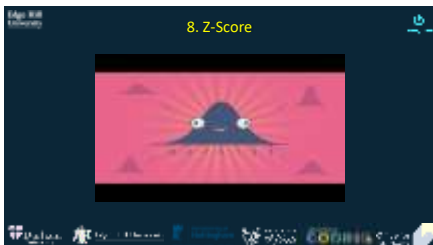
---

---

---

---

---




---

---

---

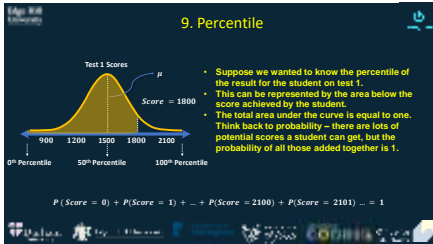
---

---

---

---

---




---

---

---

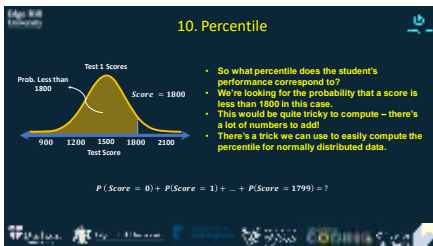
---

---

---

---

---




---

---

---

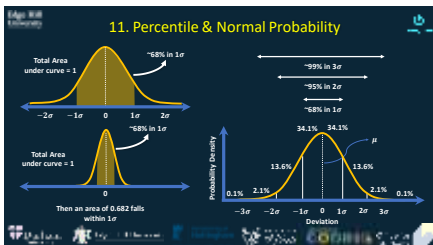
---

---

---

---

---




---

---

---

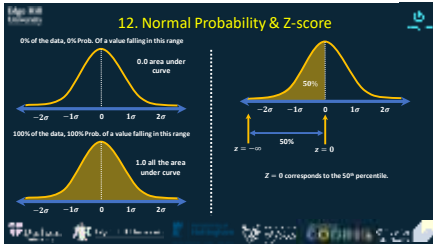
---

---

---

---

---




---

---

---

---

---

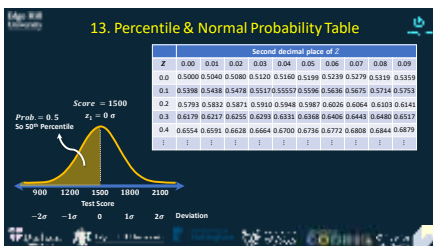
---

---

---

---

---




---

---

---

---

---

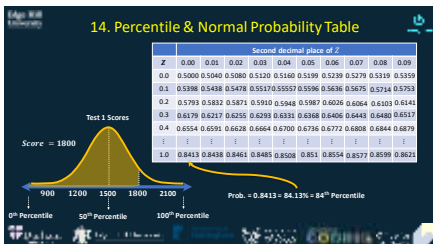
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

### 15. Percentile & Normal Probability Table

- There are two normal probability tables: for when  $z$  is negative, when  $z$  is positive.
- You don't need to remember normal probability tables.
- We can create them in code.
- What matters is that you understand that:
  - normal probability tables exist.
  - they can be used to determine what percentile an observation is in.
  - you must usually compute the  $z$ -score to make use of them.

---

---

---

---

---

---

---

---

### 16. Percentile & Normal Probability Table

- Sometimes we may not be looking for simple percentiles for our data.
- We may wish to know what proportion of our data sits between two specific positions.
- We can use the concepts we've already learned to answer some questions.
- We can do this by first calculating percentiles and then subtracting them from 1.
- Once we determine the remainder, we can use this in further calculations.

---

---

---

---

---

---

---

---

### 17. Standard Error

Standard Error of the Sample mean  $\bar{x}$

$$SE_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

Samples in observation  $\nearrow$  Standard deviation

- When we collect data, it usually represents a sample from a much larger population.
- Are our summary statistics accurate?
- The sample mean  $\bar{x}$  won't be exactly equal to the population mean  $\mu$ . It might vary from the true population quite a lot, if the sample is small.
- The standard deviation associated with an estimate is called the Standard error of an estimate.
- The standard error for  $\bar{x}$  is an important statistic – provides an indication of how uncertain we are in  $\bar{x}$ .

---

---

---

---

---

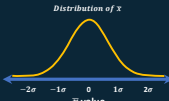
---

---

---

### 18. Confidence Intervals

- The sample mean for a collection of observations, represents an estimate of  $\mu$ .
- If we were to make another random sampling, we'll get a slightly different mean estimate.
- If we were to take many random samples from the population, and compute the sample mean for each - we would obtain a distribution for the sample mean.
- The average of this distribution is going to be very close to the true mean.
- But how confident are we in our sample mean estimate?




---

---

---

---

---

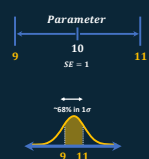
---

---

---

### 19. Confidence Intervals

- We can apply what we call "confidence intervals" to our estimates, to quantify our confidence level.
- A confidence interval contains the plausible range of values for an estimated parameter, when taking uncertainty into account using the standard error.
- For example, suppose we have an estimate for some parameter equal to 10.
- Suppose we also know the estimated parameter has a standard error of 1.
- This means it can plausibly deviate by 1.
- We can take this into account by creating an interval, that takes this deviation into account.
- The plausible range is given by the parameter plus 1, and minus 1 ( $\pm$ ).
- This is a confidence interval.




---

---

---

---

---

---

---

---

### 20. 95% Confidence Interval

Estimate  $\pm 2 \times SE_{\bar{x}}$

Parameter being estimated      Plus-minus Symbol

- We can construct a 95% confidence interval over the parameter we wish to estimate, in this case the sample mean, via the following simple formula:

$$\text{Lower Limit} \left| \frac{\text{Parameter}}{-2 \times SE_{\bar{x}} \quad + 2 \times SE_{\bar{x}}} \right| \text{Upper Limit}$$

Standard Error of the Sample mean  $\bar{x}$

$$SE_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

Samples in observation      Standard deviation

---

---

---

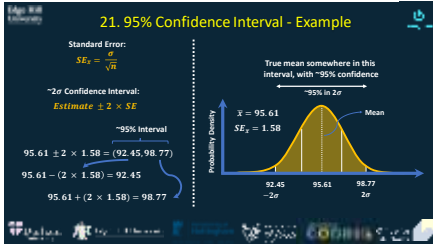
---

---

---

---

---




---

---

---

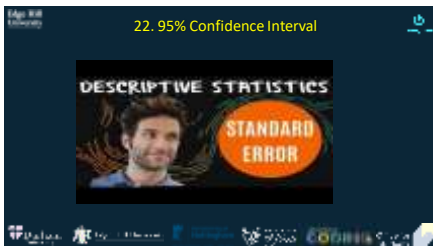
---

---

---

---

---




---

---

---

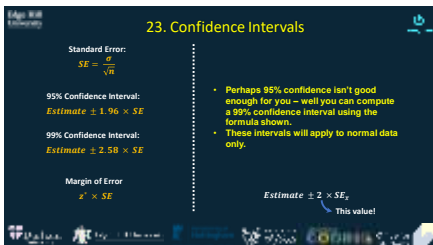
---

---

---

---

---




---

---

---

---

---

---

---

---



### 24. Testing Hypotheses

- We can start testing competing hypotheses using confidence intervals. Suppose we have a dataset describing the finishing times of runners in a race.
- We want to determine if the runners finished in a faster time this year, compared to last year.
- We form two competing hypotheses for this data. The null hypothesis is that there is no difference to average finishing times. The alternative hypothesis, is that the average runtime was different this year compared to last.
- The average runtime for last year's run was 93.29 minutes, (93 minutes and 17 seconds). We thus reframe our hypotheses given this data.

Null Hypothesis	Alternative Hypothesis
$H_0 = \text{No difference}$	$H_a \text{ or } H_1 = \text{A difference}$
$H_0: \mu_{2019} = 93.29$	$H_1: \mu_{2019} \neq 93.29$

---

---

---

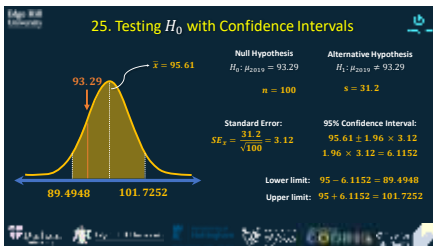
---

---

---

---

---




---

---

---

---

---

---

---

---

### 26. Decision Errors

- In general for any hypothesis test, there are four potential test outcomes.
- When running hypothesis tests, we aim to minimize the errors we make. Confidence intervals are great, but alone they don't really help us achieve that. Instead we try to use significance levels to determine how significant a result is, before making a decision.

	Do not reject $H_0$	Reject $H_0$ , Accept $H_1$
Ground Truth	$H_0$ True: Success	$H_0$ True: Type I Error
	$H_1$ True: Type II Error	$H_1$ True: Success

---

---

---

---

---

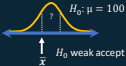
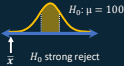
---

---

---

## 27. Decision Errors

- Confidence intervals are simplistic when it comes to hypothesis testing.
- Suppose we use a 95% confidence interval for some sample mean data, where the null hypothesis is accepted if the sample mean falls within 1 standard deviation of the mean.
- Sometimes the evidence against the null hypothesis may be overwhelming, like here.
- But sometimes we may be on the cusp of rejecting the null hypothesis, but don't quite have enough evidence to reject it.
- In these situations it's helpful to be able to quantify our confidence in the decisions we make. We can do this using a tool called, the P-value.



## 28. Decision Errors

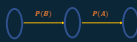


## 29. P-values

- P-values allow us to test the strength of the evidence against the null hypothesis.
- The P-value is a conditional probability – it is the probability of observing data at least as favourable to the alternative hypothesis as our current dataset is, if the null hypothesis is true.
- It may help to think of this description as a tree diagram. We can see here that the p-value is simply assessing the probability of seeing data this favourable to the alternative hypothesis, given that the null hypothesis is true.
- We usually use a summary statistic such as the sample mean to help compute a P-value.

Conditional Probability

$$P(A|B) = \frac{P(A \wedge B)}{P(B)}$$

If A = Data favourable to  $H_1$ If B =  $H_0$

### 30. P-values + Significance Level

- Example: A national sleep study suggests students sleep on average 7 hours per night.
- You're a data scientist at a local education authority, and are tasked with determining if student in your area are similar.
- You collect data from a student sample ( $n = 110$ ), and find that students in your area are sleeping on average, over seven hours.
- You want to verify that your students are indeed different from the national sample.
- You form two hypotheses.

Null Hypothesis

$H_0$ : No difference

$H_0: \mu = 93.29$

Alternative Hypothesis

$H_a$  or  $H_1$ : A difference

$H_1: \mu > 7$

---

---

---

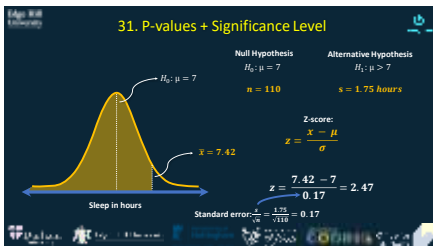
---

---

---

---

---




---

---

---

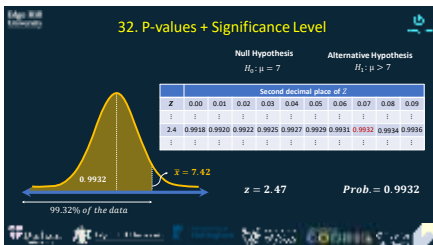
---

---

---

---

---




---

---

---

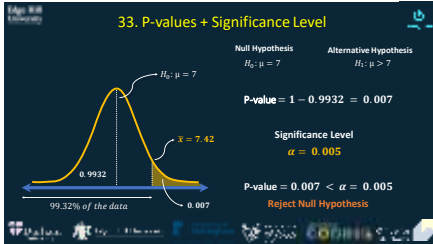
---

---

---

---

---




---

---

---

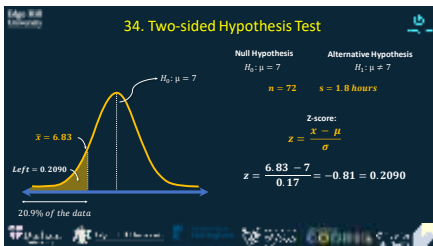
---

---

---

---

---




---

---

---

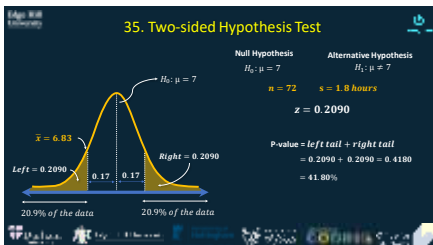
---

---

---

---

---




---

---

---

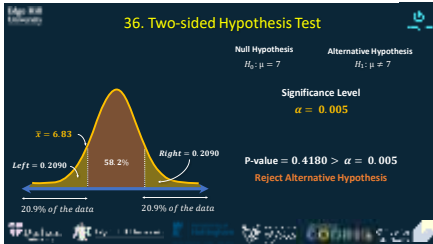
---

---

---

---

---




---

---

---

---

---

---

---

---

37. Hypothesis Testing Steps
- Some of what we've covered here may not make sense - yet. That's ok, because nobody becomes a hypothesis testing expert over night!
  - What matters is that you appreciate what's happening and why.
    - We form hypotheses to answer questions about our data.
    - We collect data samples to test them.
    - We compute summary statistics over the data sample, such as the sample mean and sample standard deviation.
    - We compute the Z-score and use this along with normal probability tables to determine the area under the curve.
    - We use these areas to represent probabilities as p-values, and evaluate them with respect to some significance level,  $\alpha$ .
  - A little practice will help make these ideas clearer.

---

---

---

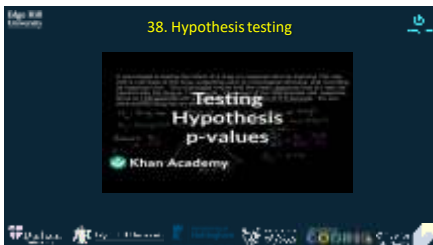
---

---

---

---

---




---

---

---

---

---

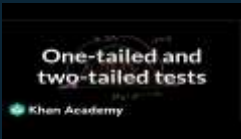
---

---

---

Edg 9.0  
University

### 39. Hypothesis testing



One-tailed and two-tailed tests

Khan Academy

Windows taskbar icons: Edge, File Explorer, Microsoft Store, etc.

---

---

---

---

---

---

---

---

Edg 9.0  
University

### 40. Activities



Link to the notebook:

Windows taskbar icons: Edge, File Explorer, Microsoft Store, etc.

---

---

---

---

---

---

---

---

Edg 9.0  
University


### 41. Resources

**Books:**

- "OpenIntro Statistics, 4th ed", D. Olney, M. Golinberg-Romel and C. Burt.
- "Data Science From Scratch: First Principles with Python", 2nd Edition, J. Gras.
- "Think Stats: Probability and Statistics for Programmers", A. & Downey.
- "Statistics in Plain English, Third Edition, Volume 17", T. C. Urdan.

**Tools Websites**

- [Kaggle](#) - an online platform where you can tackle data science challenges.
- [Kward data science](#) - a website where data science practitioners share ideas, tutorials and advice.



Windows taskbar icons: Edge, File Explorer, Microsoft Store, etc.

---

---

---

---

---

---

---

---

Edg 9.0  
University

## 42. Checkpoint

We've reached another checkpoint. Let's recap what we've introduced so far.

- Normal distributions.
- The  $Z$  score.
- Probability tables.
- Standard error.
- Confidence intervals.
- Hypothesis testing.

From here you can pursue the activities provided in Google Colab, or watch the next set of slides which cover the ethics of data science. It's entirely up to you.

UCLouvain | UCLouvain | UCLouvain | UCLouvain | UCLouvain | UCLouvain | UCLouvain | UCLouvain | UCLouvain | UCLouvain

---

---

---

---

---

---

---