



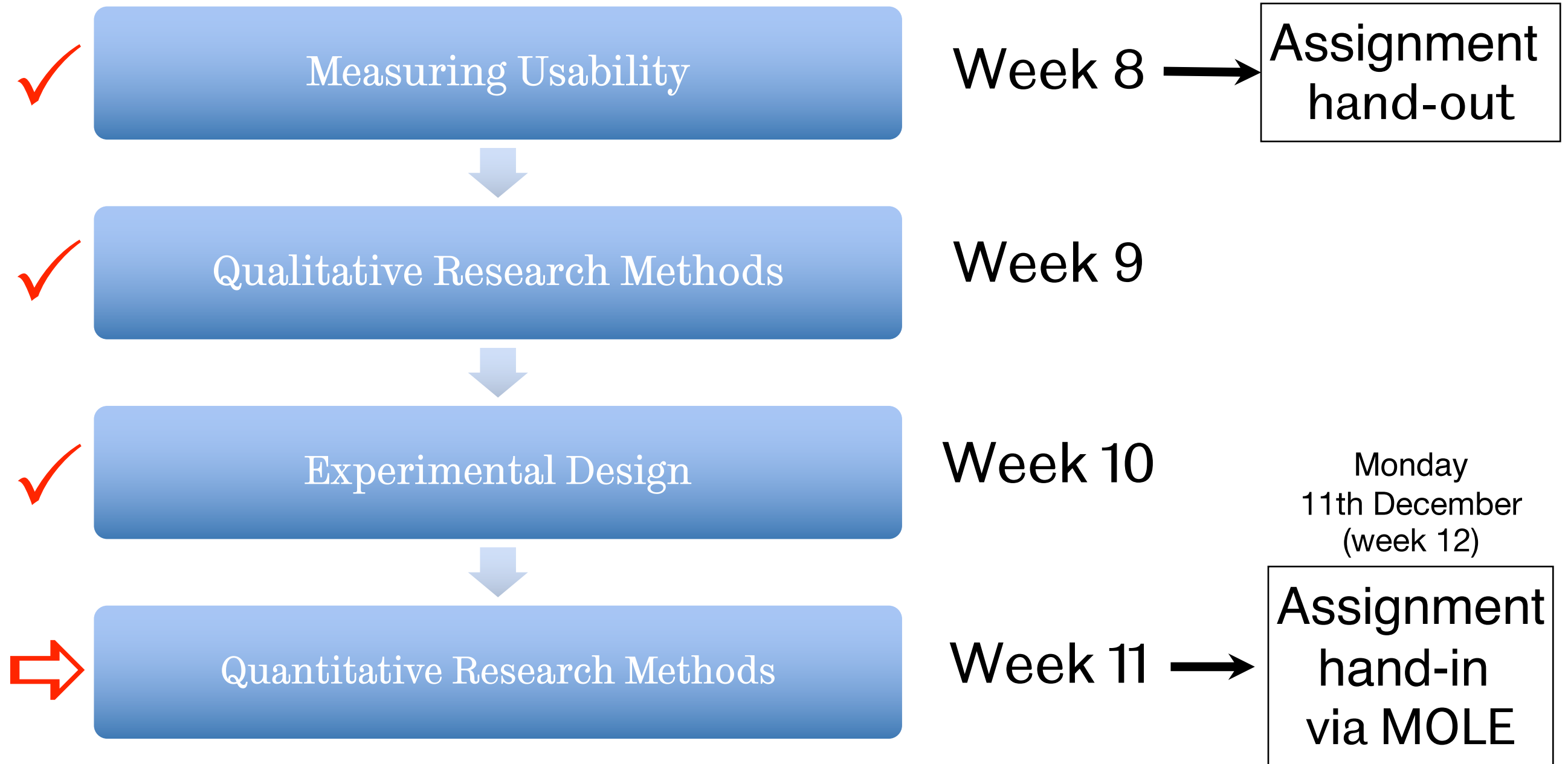
The  
University  
Of  
Sheffield.

# Human Centred Systems Design

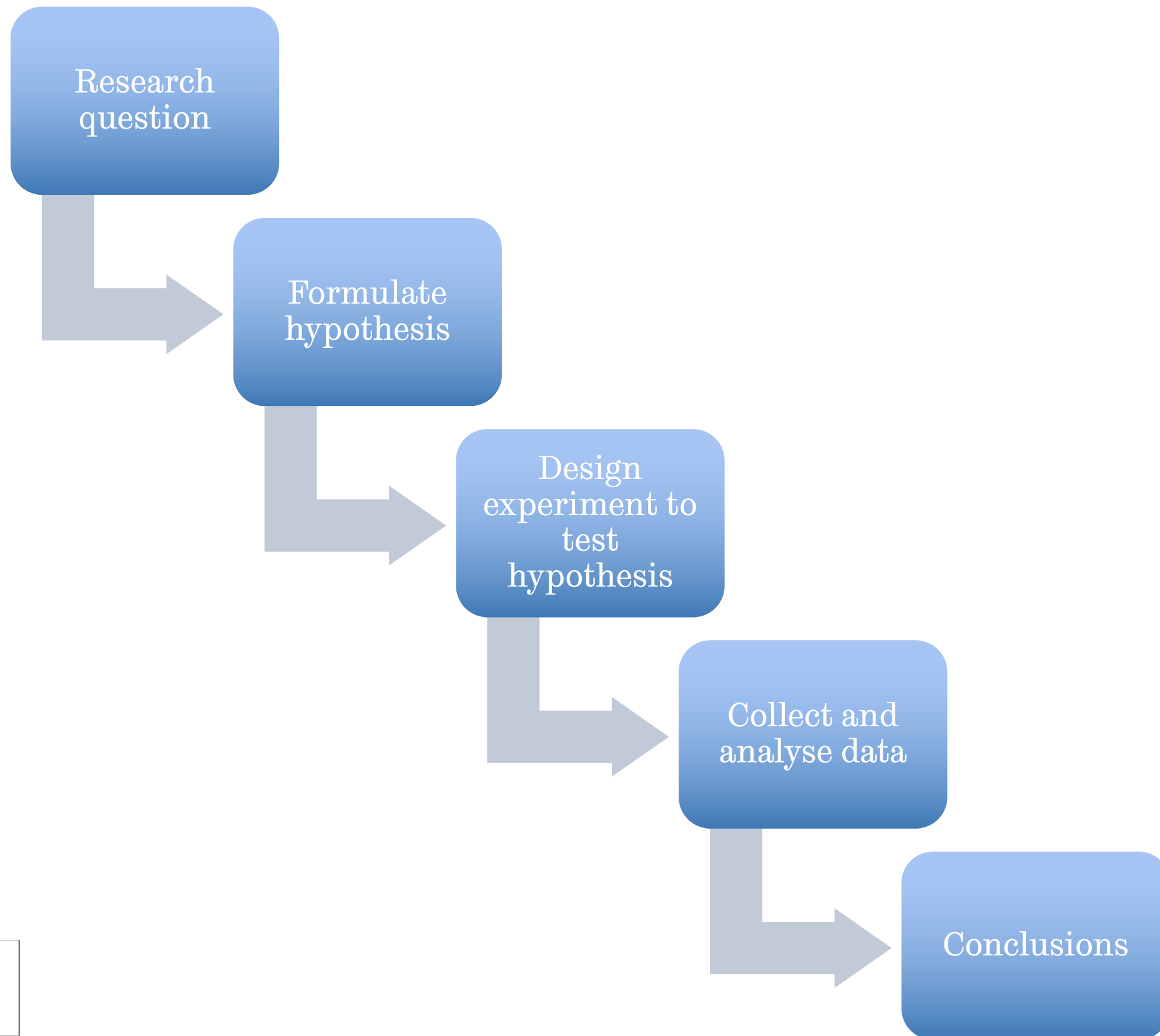
## Quantitative Research Methods

Dr Maria-Cruz Villa-Uriol

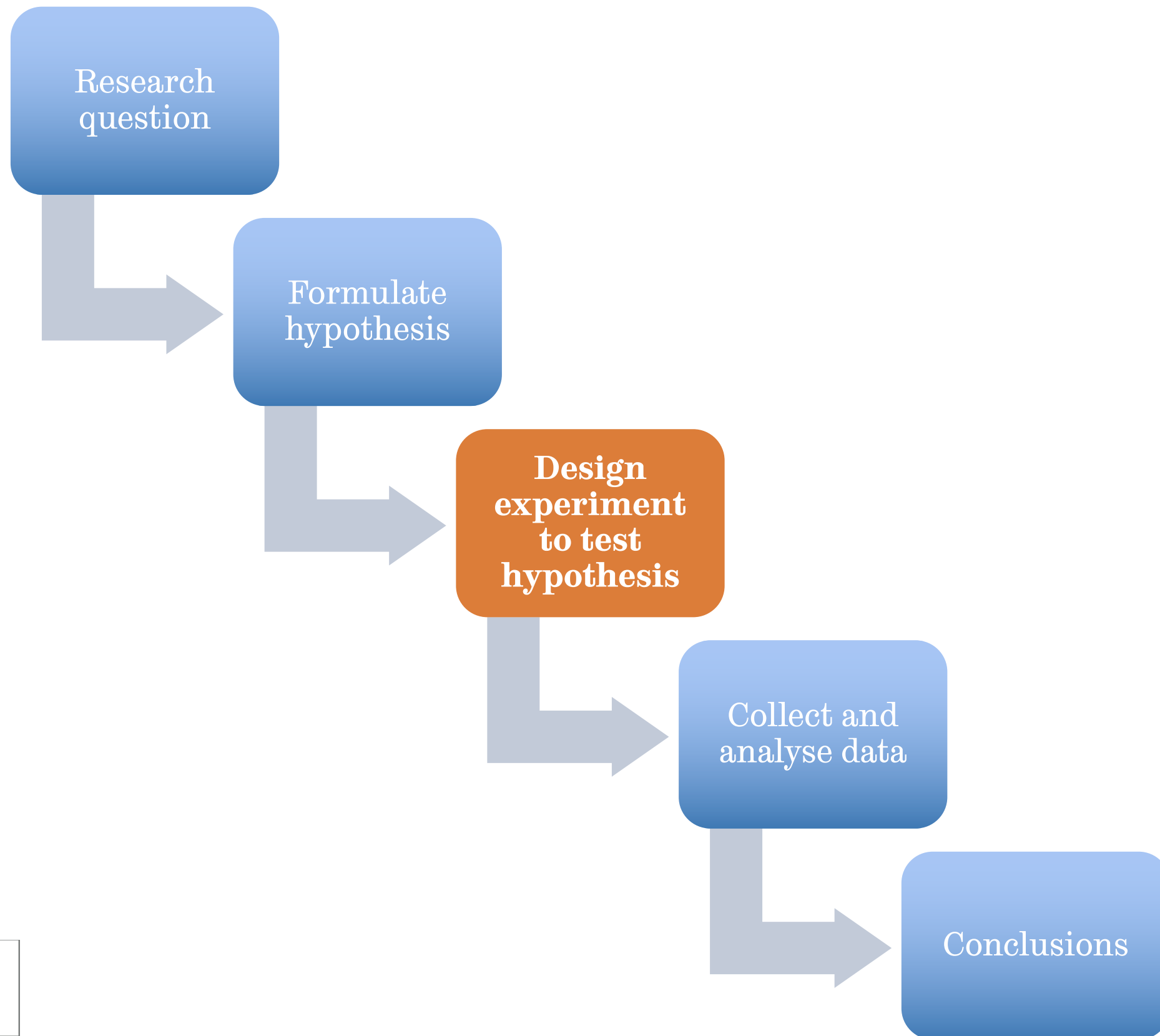
# Evaluation in HCI



# Experimental Research Methods



# Experimental Research Methods



Qualitative research	vs.	Quantitative research
Discover ideas and gain insight and understanding	<b>Aim</b>	Test hypotheses and specific research questions
Observe, survey and interpret	<b>Approach</b>	Measure and test
Mixed	<b>Data collection</b>	Structured
Researcher involved and results subjective	<b>Researcher independence</b>	Researcher uninvolved observer, objective results
Small samples, naturalistic setting	<b>Sample size</b>	Larger samples for generalisable results

# Qualitative Research Methods

- Qualitative vs Quantitative Data
- Qualitative Research Methods
- Traditionally, in Human-Computer Interaction:
  - ① Surveys
  - ② Interviews
  - ③ Focus groups
  - ④ Diaries
  - ⑤ Ethnographic research

**Mostly descriptive!**

# Quantitative Research Methods

1

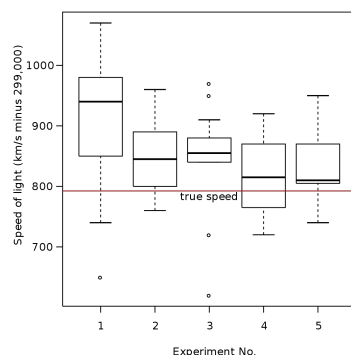
## • Descriptive statistics

- organising data
- summarising data
- simplifying data
- describing and presenting data

2

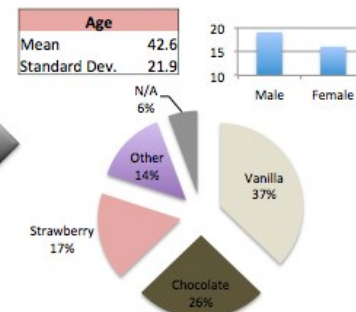
## • Inferential statistics

- generalising from samples to populations
- making predictions
- hypothesis testing



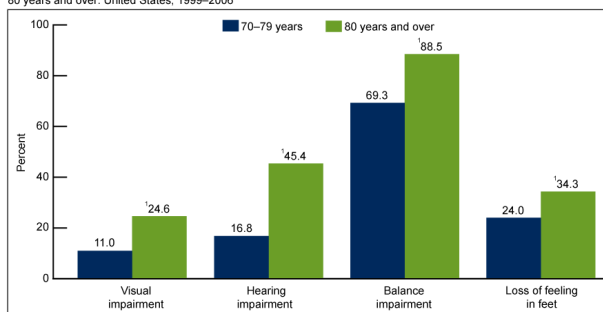
Respondent #	Age	Gender	Favorite Ice Cream Flavor
1	36	m	Vanilla
2	22	f	Chocolate
3	61	m	Strawberry
4	88	m	Other
5	31	m	N/A
6	53	m	N/A
7	30	f	Chocolate
8	64	f	Chocolate
9	18	m	Vanilla
10	16	f	Vanilla
11	83	m	Strawberry
12	16	f	Strawberry
13	94	m	Strawberry
14	55	m	Vanilla
15	42	f	Chocolate
16	18	f	Vanilla
17	19	f	Vanilla

Raw Data

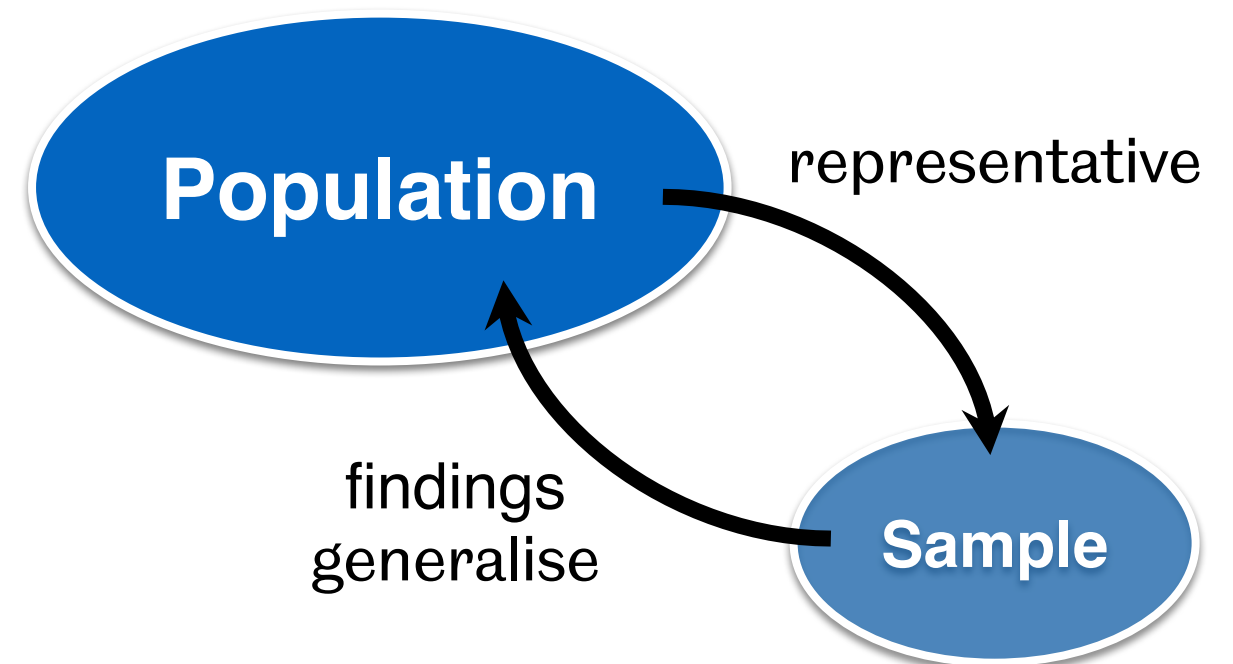
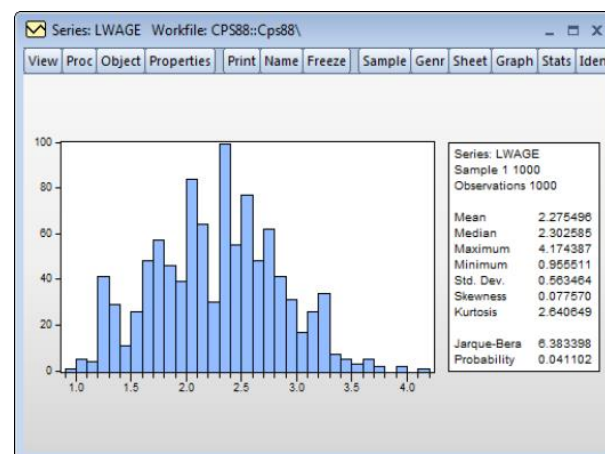


Descriptive Statistics

Figure 3. The prevalence of sensory impairments among persons aged 70–79 years compared with persons aged 80 years and over: United States, 1999–2006

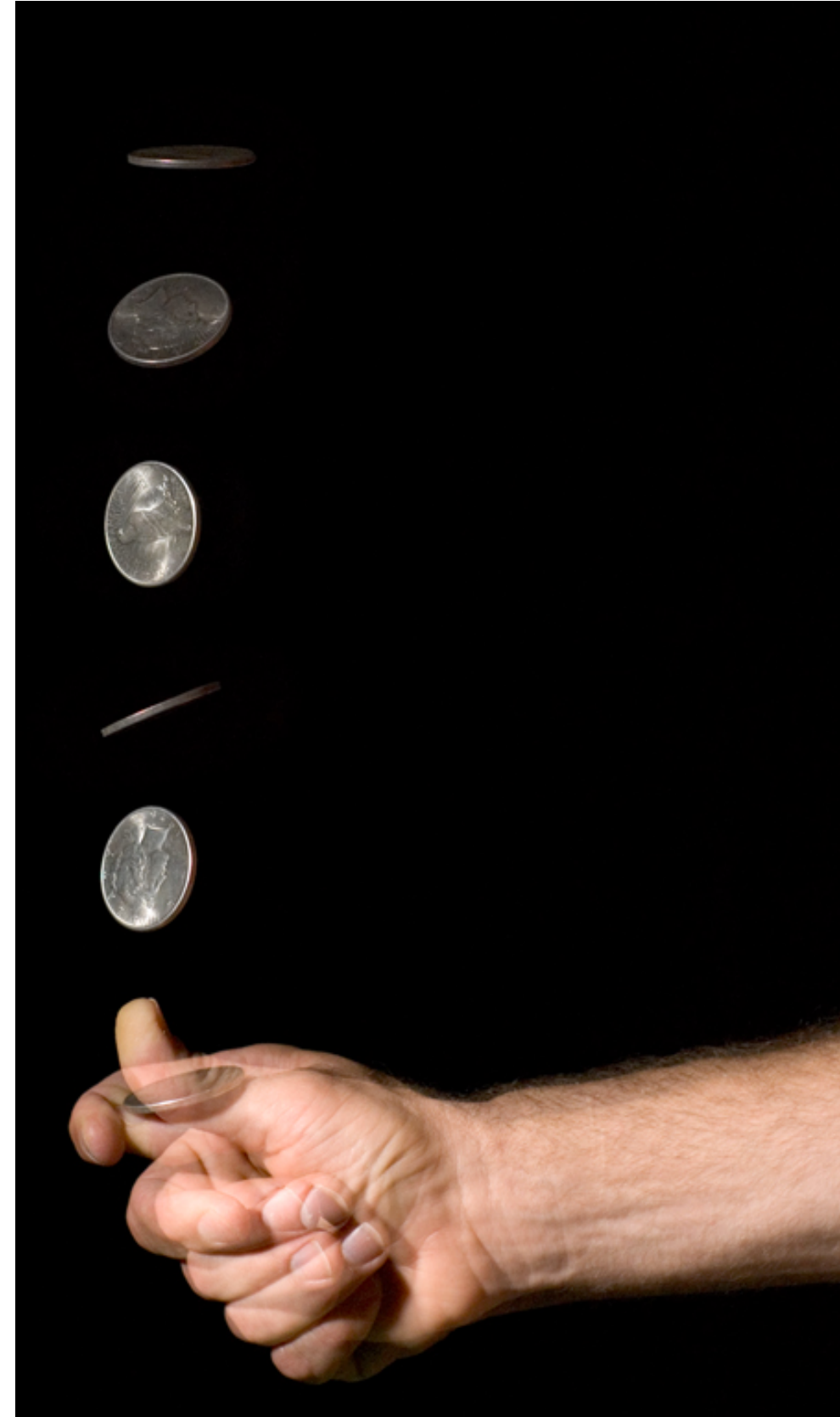


\*Significantly different from the 70–79 age group.  
SOURCE: CDC/NCHS, National Health and Nutrition Examination Survey.



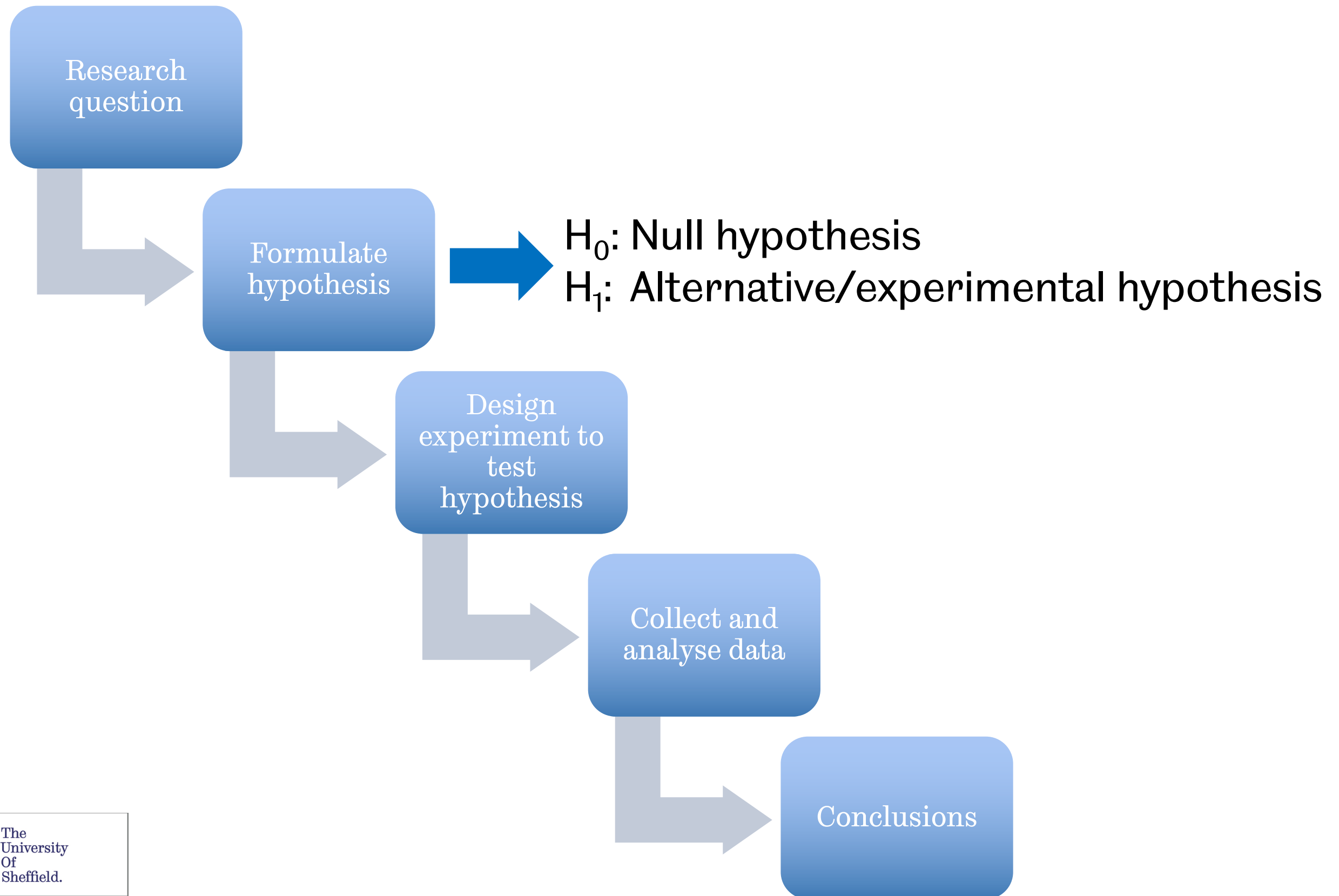
# Example: Tossing a coin

- Toss a coin 10 times and record the order of heads (H) and tails (T)
- Which is more likely?
  - HHHHHHHHHH?
  - TTTTTTTTTT?
  - HTHTHTHTHT?
- Research question:  
Is this coin fair or in favour of heads?





# Experimental Research Methods



# Null hypothesis

- $H_0$  – The hypothesis that nothing has changed
  - E.g. That our actions have not had a corresponding effect
  - E.g. There is no difference between experimental treatments

We can never prove the null hypothesis  
We may just find evidence to reject it

# Experimental (alternative) hypothesis

- $H_1$  – The hypothesis we are interested in, what we are testing for
  - E.g. That something has changed in the world because of our actions

It is mutually exclusive with the null hypothesis,  
i.e. they can not happen simultaneously

# Null hypothesis: Example

- Research question: do dogs eat bananas?
- $H_1$  – Experimental/alternative hypothesis: dogs eat bananas
- $H_0$  – Null hypothesis: dogs do not eat bananas
- Null hypothesis is tested –  
It can be rejected if a dog is seen eating a banana



# Null hypothesis: Example

- We believe that people will be able to remember more words from a list when (a) they are organised into meaningful groups than when (b) they are presented randomly.
- H<sub>0</sub> – Null Hypothesis:
  - There is no difference in the number of words people can remember under condition (a) and condition (b)
- H<sub>1</sub> – Experimental/Alternative Hypothesis:
  - People will remember more words in condition (a) than in condition (b)

# Example: Tossing a coin

- Research question:

Is this coin fair or biased in favour of heads?

- $H_0$  - null hypothesis: the coin is fair

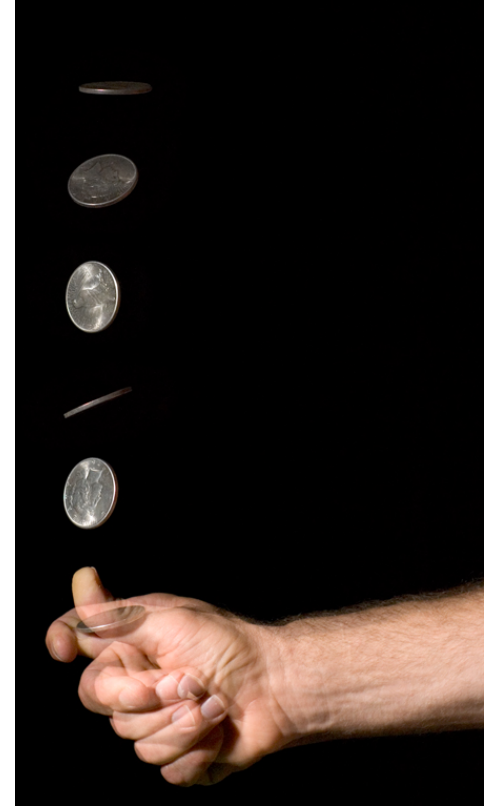
(i.e., the probability of a head is 0.5)

- $H_1$ : alternative/experimental hypothesis: the coin is

biased in favour of a head (i.e. the probability of a head is greater than 0.5, e.g. it is 0.7)

- Experiment: Toss a coin 10 times and record the order of heads (H) or tails (T)

- Testing: On the basis of the outcomes, make your decision.



# Example: Tossing a coin

- Research question:

**But...**

**How many times should we  
toss the coin to be sure?  
Is 10 times sufficient?**

order of heads (H) or tails (T)

- Testing: On the basis of the outcomes, make your decision.

# Example (and a Task)

a) Mike's height is 6'2". Mary's height is 5'8".

***So Mike is taller than Mary.***

b) The average height of three males (Mike, John, and Ted) is 5'5". The average height of three females (Mary, Rose, and Jessica) is 5'10".

***So females are taller than males.***

Some questions:

- *Which of the above statements are correct?*
- *Why?*



# Example (possible answers)

b) The average height of three males (Mike, John, and Ted) is 5'5". The average height of three females (Mary, Rose, and Jessica) is 5'10".  
*So females are taller than males.*

- It's common sense, males are generally taller than females
- I can easily find 3 other males and 3 other females, in which the average height of the males is higher than that of the females
- There are only 3 individuals in each group! The sizes of the comparison groups are too small
- The individuals in both the male and female groups are not representative of the general population

***Clear example of inappropriate sampling!***

# Why would I want to do this in HCI?

- Because we want to demonstrate something
- e.g.
  - Improvement in processing
  - Better user interface design
  - More acceptable robot
- We want robust findings that allow us to draw reliable conclusions

# Significance Tests

- Purpose:

Indicate whether observed differences between assessment results occur because of sampling error or chance.

- When are they necessary?

- When all values of the members of the comparison groups are known, direct comparison of values is possible and we can draw a conclusion.

*Not needed, there is no uncertainty involved*

- When the population is large, we can only sample a sub-group of people from the entire population

*Needed, as they allow us to determine how confident we are that the results observed from the sampling population can be generalized to the entire population*

# Example (and a Task)

- a) Mike's height is 6'2". Mary's height is 5'8".
- b) The average height of three males (Mike, John, and Ted) is 5'5".  
The average height of three females (Mary, Rose, and Jessica) is 5'10".

- When all values of the members of the comparison groups are known, direct comparison of values is possible and we can draw a conclusion.

*Not needed, there is no uncertainty involved*

**a) or b)????**

- When the population is large, we can only sample a sub-group of people from the entire population

*Needed, as they allow us to determine how confident we are that the results observed from the sampling population can be generalized to the entire population*

**a) or b)????**

# Example (and a Task)

- a) Mike's height is 6'2". Mary's height is 5'8".
- b) The average height of three males (Mike, John, and Ted) is 5'5".  
The average height of three females (Mary, Rose, and Jessica) is 5'10".

- When all values of the members of the comparison groups are known, direct comparison of values is possible and we can draw a conclusion.

*Not needed, there is no uncertainty involved*

- a) Mike's height is 6'2". Mary's height is 5'8".

- When the population is large, we can only sample a sub-group of people from the entire population

*Needed, as they allow us to determine how confident we are that the results observed from the sampling population can be generalized to the entire population*

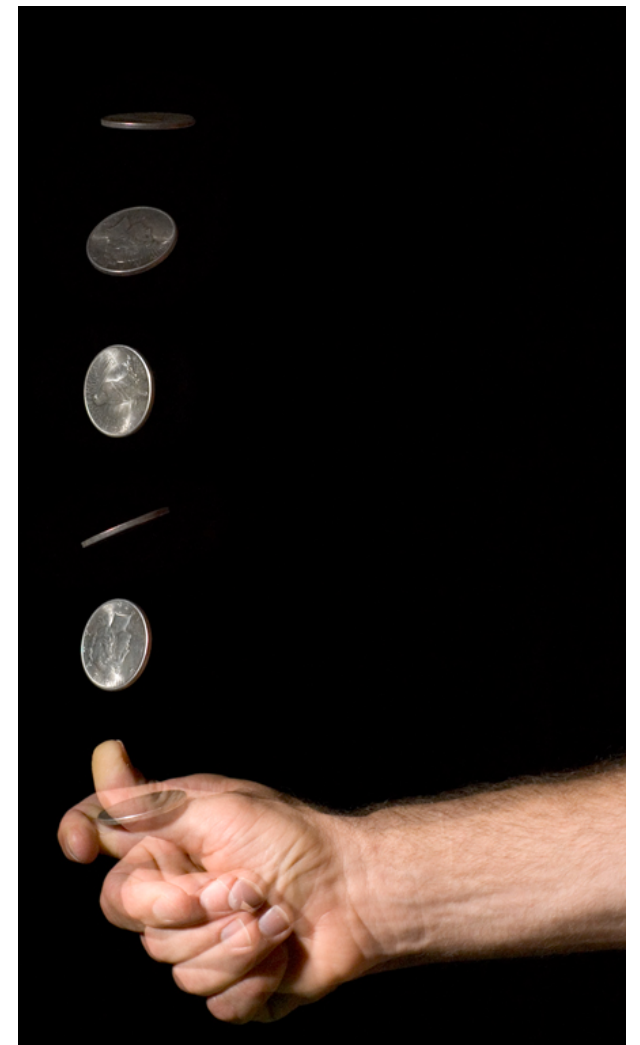
- b) The average height of three males (Mike, John, and Ted) is 5'5".  
The average height of three females (Mary, Rose, and Jessica) is 5'10".

# Example: Tossing a coin

- Research question:

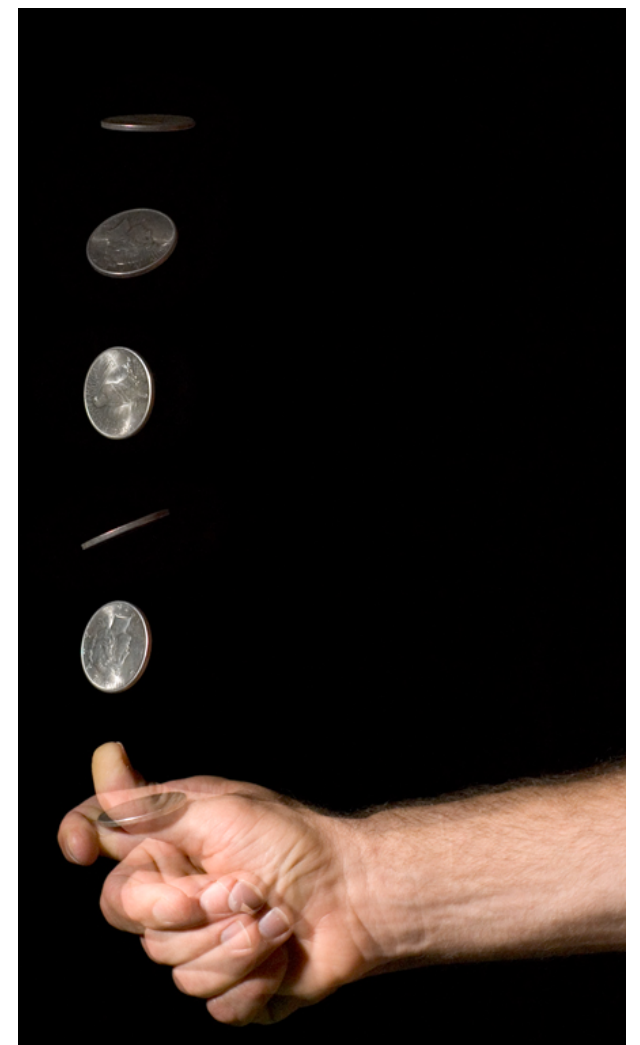
Is this coin fair or biased in favour of heads?

- Experiment: Toss a coin 10 times and record the order of heads (H) or tails (T)



# Tossing a coin: Possible outcomes

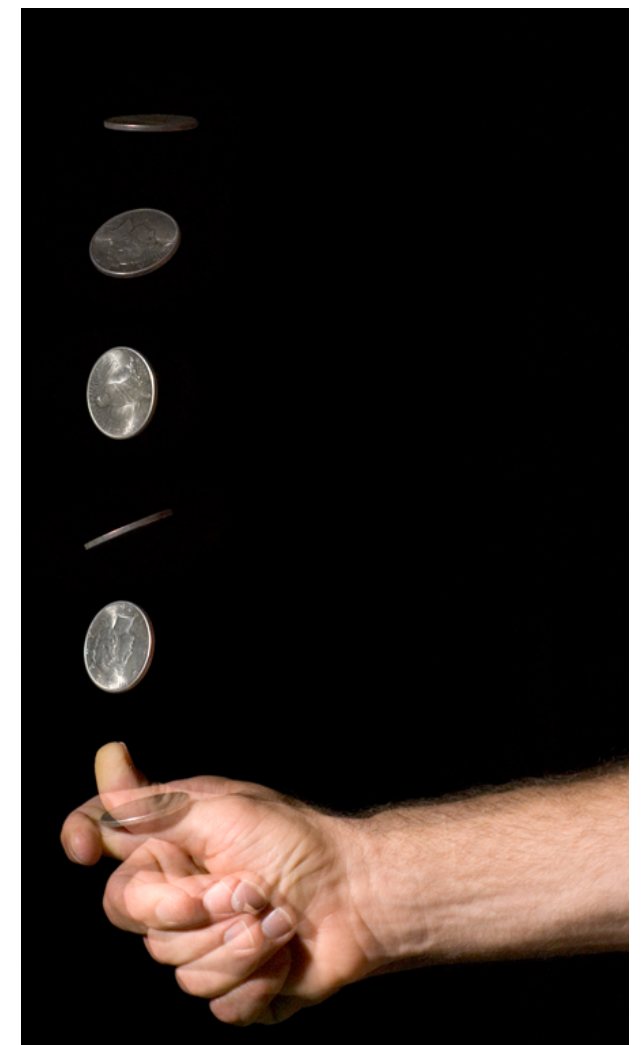
- Coin is not biased
  - Our data indicates that it is not biased
    - Correct decision
  - Our data indicates that it is biased
    - Incorrect decision
- The incorrect case is an example of Type I error (false positive)





# Tossing a coin: Possible outcomes

- Coin is biased
  - Our data indicates that it is biased
    - Correct Decision
  - Our data indicates that it is not biased
    - Incorrect Decision
- The incorrect case is an example of Type II Error (false negative)





# Type I and Type II errors

- All significance tests are subject to the risk of Type I and Type II errors
  - Type I error ( $\alpha$  error or “false positive”):  
*The mistake of rejecting the null hypothesis when it is true and should have not been rejected*
  - Type II error ( $\beta$  error or a “false negative”):  
*The mistake of not rejecting the null hypothesis when it is false and should have been rejected*

# Type I and Type II errors: example

- $H_0$ : The defendant is innocent
- $H_1$ : The defendant is guilty

A judicial case (innocence is presumed)		Jury decision ("prediction")	
		Not guilty	Guilty
Reality (truth)	Not guilty	✓	Type I error
	Guilty	Type II error	✓

# Type I and Type II errors: example

- $H_0$ : No difference between ease of use of ATMS with touch screens vs ATMs with buttons
- $H_1$ : ATMs with touch screens are easier to use than ATMs with buttons

HCI experiment

Reality (truth)		Study conclusion ("prediction")	
		???	???
		✓	Type I error
	???	Type II error	✓

# Type I and Type II errors: example

- $H_0$ : No difference between ease of use of ATMS with touch screens vs ATMs with buttons
- $H_1$ : ATMs with touch screens are easier to use than ATMs with buttons

HCI experiment

		Study conclusion ("prediction")	
		No difference	Touchscreen ATM easier
Reality (truth)	No difference	✓	Type I error
	Touchscreen ATM easier	Type II error	✓

# Type I and Type II errors: example

- $H_0$ : Patient does not suffer from disease A
- $H_1$ : Patient does suffer from disease A

Diagnosis test		Outcome of diagnosis test ("prediction")	
		Negative	Positive
Reality (truth)	Negative	ppp	ppp
	Positive	ppp	ppp

# Type I and Type II errors: example

- $H_0$ : Patient does not suffer from disease A
- $H_1$ : Patient does suffer from disease A

Diagnosis test		Outcome of diagnosis test ("prediction")	
		Negative	Positive
Reality (truth)	Negative	✓	Type I error
	Positive	Type II error	✓



# Type I and Type II errors

- Typical considerations:
  - It is generally believed that Type I errors are worse than Type II errors
  - Statisticians call Type I errors a mistake that involves “gullibility”
    - A Type I error may result in a condition worse than the current state
  - Type II errors are mistakes that involve “blindness”
    - A Type II error can cost the opportunity to improve the current state

# Controlling risks of errors

- In statistics,
  - the probability of making a Type I error is called alpha ( $\alpha$ ) or significance level ( $p$ -value)
  - the probability of making a Type II error is called beta ( $\beta$ )
  - the statistical power of a test, defined as  $1-\beta$ , refers to the probability of successfully rejecting a null hypothesis when it is false and should be rejected

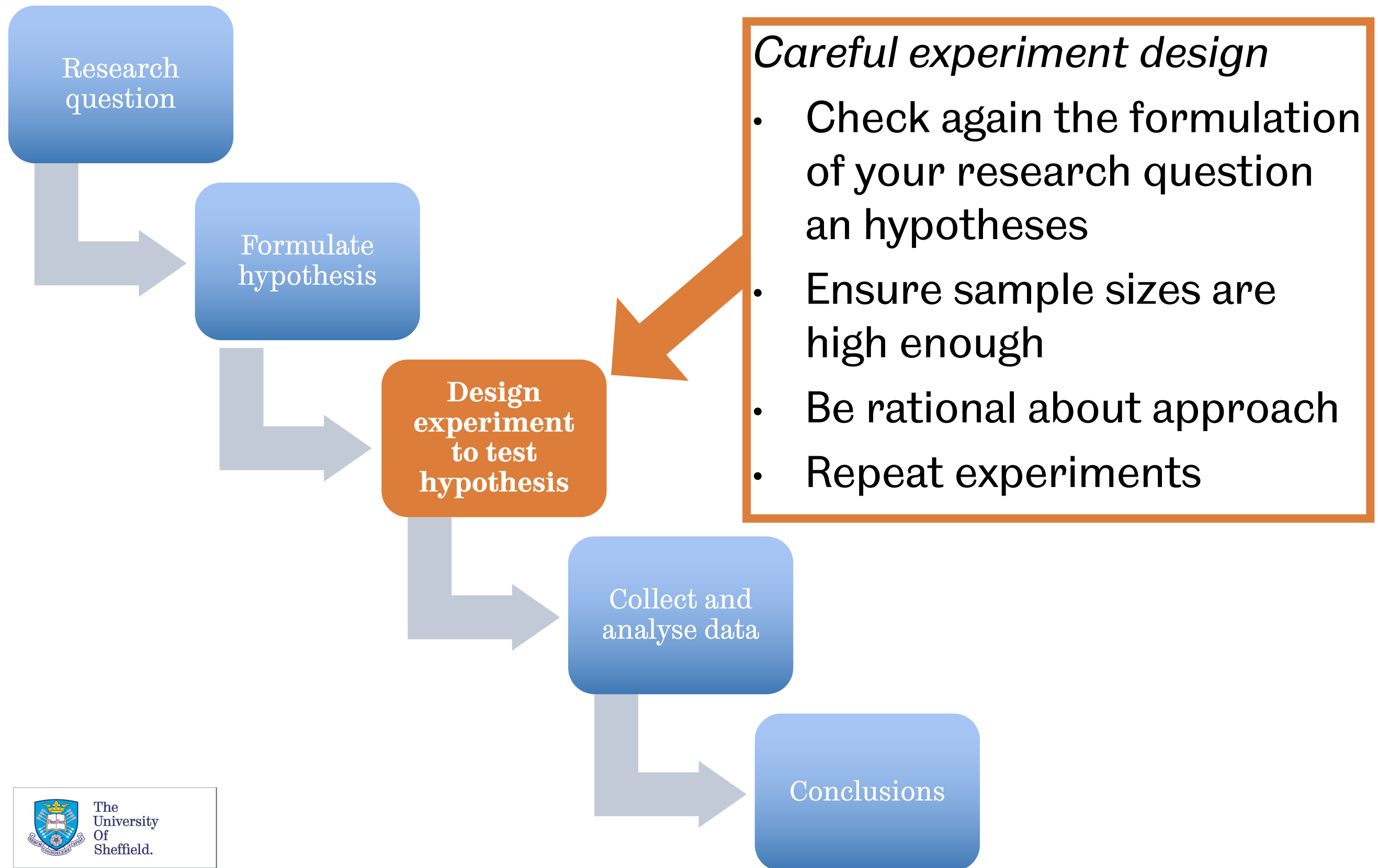


# Controlling risks of errors

- Alpha ( $\alpha$ ) and beta ( $\beta$ ) are interrelated and under the same conditions:
  - a decrease in  $\alpha$  reduces the chance of making Type I errors,
  - but increases the chance of making Type II errors
- In experimental research, it is generally believed that
  - Type I errors are worse than Type II errors

So a very low  $p$ -value (0.05) is widely adopted to control the occurrence of Type I errors

# Controlling risks of errors



# Example: Research Question

The developers of a website want to know if a pull-down menu or a pop-up menu is better

- Research question:
  - Which menu type is best for navigating the site?



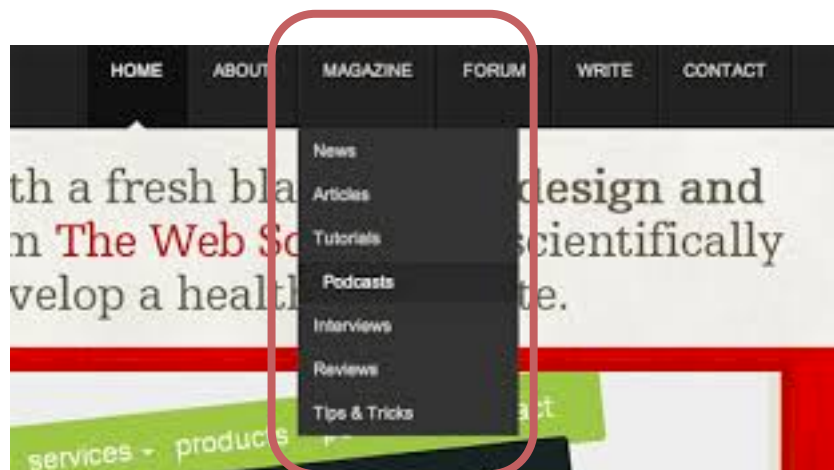
VS.



# Formulating a hypothesis: Example

The developers of a website want to know if a pull-down menu or a pop-up menu is better

- Research question:
  - Which menu type is best for navigating the site?
- What we evaluate (*independent variable*):
  - the type of menu (pull-down or pop-up)
- On the basis of what measure (*dependent variable*):
  - time spent locating web pages



VS.



# In HCI, we typically evaluate

- Technology
  - e.g. types of technology or device, types of design...
- Users:
  - e.g. age, gender, computer experience, professional domain, education, culture, motivation, mood, and disabilities
- Context of use:
  - e.g. physical status, user status, social status...

*Independent variables*

# In HCI, measures typically used

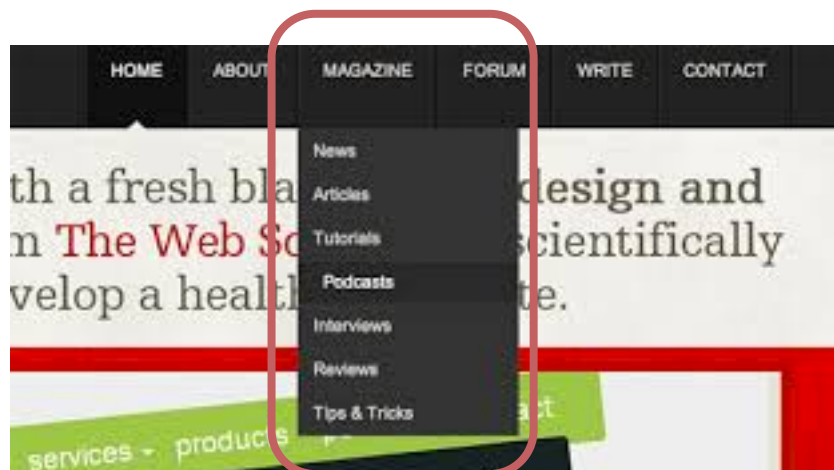
- Efficiency:
  - e.g., task completion time, speed
- Accuracy:
  - e.g., error rate
- Subjective satisfaction:
  - e.g., Likert scale ratings
- Ease of learning and retention rate
- Physical or cognitive demand
  - e.g., NASA task load index

*Dependent variables*

# Formulating a hypothesis: Example

The developers of a website want to know if a pull-down menu or a pop-up menu is better

- Research question:
  - Which menu type is best for navigating the site?
- What we evaluate (*independent variable*):
  - the type of menu (pull-down or pop-up)
- On the basis of what measure (*dependent variable*):
  - time spent locating web pages



VS.



# Example: $H_0$ and $H_1$ hypotheses

The developers of a website want to know if a pull-down menu or a pop-up menu is better

- Research question:
  - Which menu type is best for navigating the site?
- $H_0$  (null hypothesis):
  - There is no difference between both menu types in the time spent locating pages
- $H_1$  (alternative hypothesis):
  - There is a difference between both menu types in the time spent locating pages



# Example: $H_0$ and $H_1$ hypotheses

The developers of a website want to know if a pull-down menu or a pop-up menu is better

- Research question:

- Which menu type is best for navigating the

independent variable

- $H_0$  (null hypothesis):

- There is no difference between both menu types in the time spent locating pages

dependent variable

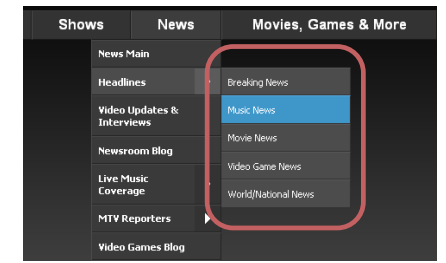
- $H_1$  (alternative hypothesis):

- There is a difference between both menu types in the time spent locating pages

# Task



vs.



The developers of a website want to know if a pull-down menu or a pop-up menu is better

- Research question:
  - Which menu type is best for navigating the site?
- What we evaluate (*independent variable*):
  - the type of menu (pull-down or pop-up)

*Please, work in pairs for 5 min and propose other metrics (dependent variables) and build the corresponding null ( $H_0$ ) and experimental/alternative hypotheses ( $H_1$ ).*

# Example: $H_0$ and $H_1$ hypotheses

The developers of a website want to know if a pull-down menu or a pop-up menu is better

- Research question:
  - Which menu type is best for navigating the site?
- $H_0$  (null hypothesis):
  - There is no difference in user satisfaction rating between the pull-down menu and the pop-up menu
- $H_1$  (alternative hypothesis):
  - There is a difference in user satisfaction between the pull-down menu and the pop-up menu

# Example: $H_0$ and $H_1$ hypotheses

The developers of a website want to know if a pull-down menu or a pop-up menu is better

- Research question:

independent variable

- Which menu type is best for navigating the site?

- $H_0$  (null hypothesis):

dependent variable

- There is no difference in user satisfaction rating between the pull-down menu and the pop-up menu

- $H_1$  (alternative hypothesis):

- There is a difference in user satisfaction between the pull-down menu and the pop-up menu

# Example: Experiment Design

- Two alternatives are possible:

## Between subjects design

1. Recruit two groups of people:
  - Group 1 will see the version of the website with pop up menus,
  - Group 2 will see a version of the website with pull down menus

## Within subjects design

2. Recruit only one group of people:
  - All participants will experience both websites

# Randomisation

- Randomisation refers to the random assignment of treatments to the experimental units or participants (Oehlert 2000)
- In a totally randomized experiment, no one, including the investigators themselves, is able to predict the condition to which a participant is going to be assigned
- Available strategies include:
  - *traditional*: tossing a coin, throwing dice, drawing capsules out of an urn...
  - *software-based*: based on the generation of random numbers
- In HCI, randomisation is typically applied in:
  - selecting participants
  - designing the tasks
  - altering the order of exposure to experiments
  - ...

# Randomisation: example

- Randomized control trial
  - participants randomly allocated to various conditions in the study
- Counter example: testing a drug vs placebo.  
Giving drug to those whose faces you liked, and placebo to those with faces you did not like wouldn't be a good method

# Significance Tests: Interpretation

- Inferential statistical tests
  - All calculate the probability of an outcome occurring.
    - Zero probability –
    - Probability of 1 –
    - Probability of .50 –



# Significance Tests: Interpretation

- Inferential statistical tests
  - All calculate the probability of an outcome occurring.
    - Zero probability – something will not happen
    - Probability of 1 – something will happen
    - Probability of .50 – happens half the time.
      - E.g. Unbiased coin, .50 probability of Heads occurring

# Significance Tests

- Statistical test will tell you the probability of seeing these scores if the null hypothesis is correct.
- E.g.  $p=0.16$ 
  - 16% chance of seeing this difference between conditions if the null hypothesis is true
  - or
  - 16% chance that we will be wrong if we conclude that the experimental hypothesis is true

# Example

Do men and women differ in their driving ability?

- Experiment:
  - Test in driving simulator:
    - 30 men make 6 errors each on average,
    - 30 women make 4 errors each on average
- Statistical test:
  - Run statistical test to compare the scores
    - $p = 0.16$

*Please, discuss in pairs for 10 min, what the null hypothesis could be, and provide a plausible interpretation of the p-value*

# Example

Do men and women differ in their driving ability?

- Experiment:
  - Test in driving simulator:
    - 30 men make 6 errors each on average,
    - 30 women make 4 errors each on average
- Statistical test:
  - Run statistical test to compare the scores
    - $p = 0.16$ 
      - 16% probability that we will be wrong if we conclude that men and women drive differently
  - or
  - 16% chance of obtaining these scores if the null hypothesis (no difference) is true

# Acceptable risk levels

- Level of risk you will accept (or cutoff for significance):
  - alpha  $\alpha$
- In human research, standard level of alpha  $\alpha$  :
  - e.g.  $\alpha = 0.05$
  - Interpretation:
    - 5% chance of conclusions being wrong
    - also expressed as a 95% confidence interval
- In medicine, level of risk accepted tends to be lower:
  - e.g.  $\alpha = 0.01$

# Levels of certainty

- Generically expressed as a p-value
  - e.g.  $p < 0.05$
- Read as:
  - The probability that, given the null hypothesis, the results seen would have arisen
- Do not read as:
  - The probability that there is an effect
- Subtle but significant distinction

# Statistical significance and p-values

- A small **p-value** (typically  $\leq \alpha$ ) **indicates** strong evidence against the null hypothesis, so you reject the null hypothesis (i.e. experiment indicates that a significant difference exists)
- A **large p-value** ( $> \alpha$ ) **indicates** weak evidence against the null hypothesis, so you fail to reject the null hypothesis
- We typically seek:
  - p value lower than  $\alpha$
- In HCI usually  $p < 0.05$ 
  - E.g. (possibly) there is a significant difference in the number of words people can remember from a random list, and the number they can remember from an organised list.

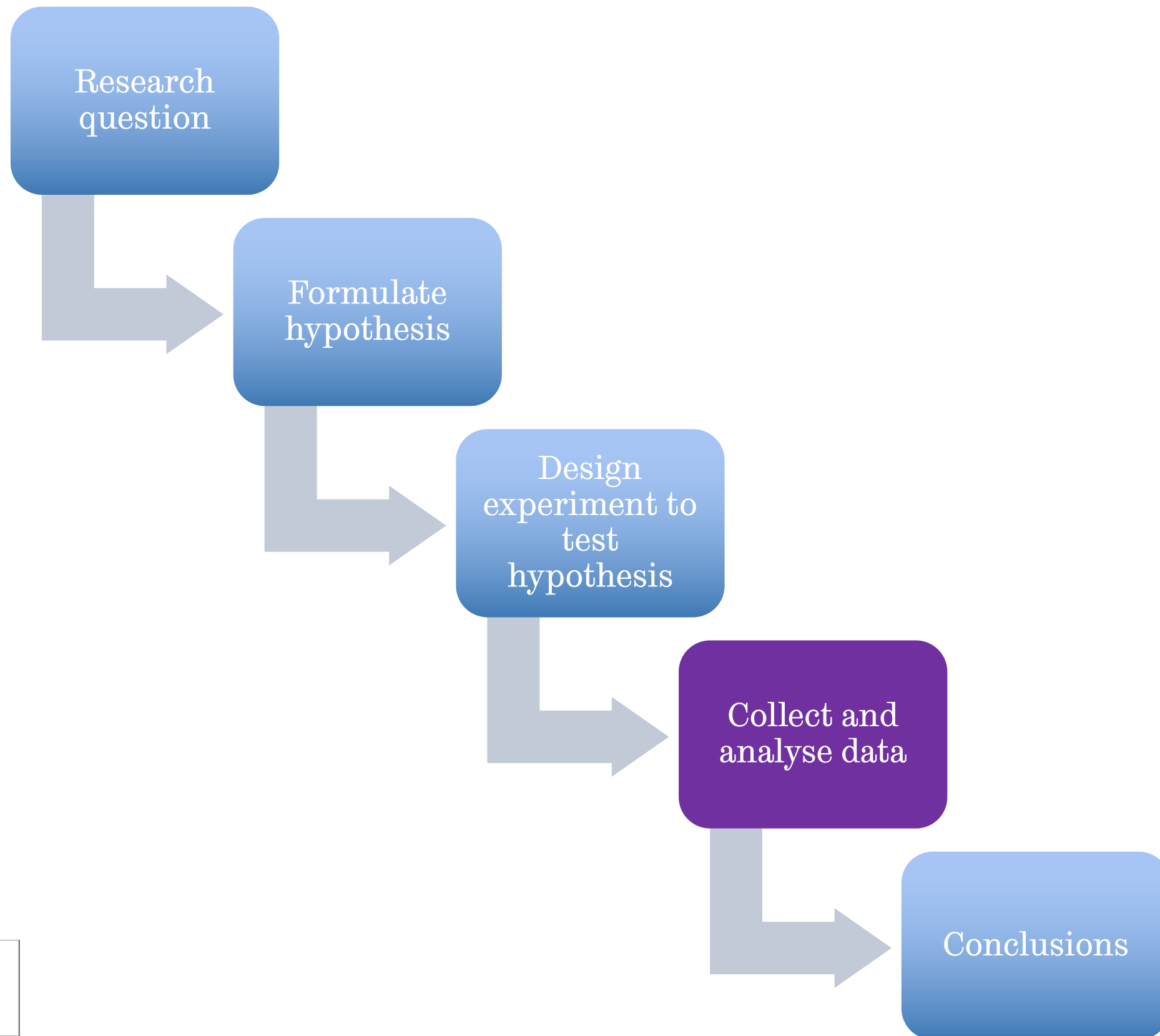
# Computation of p-values

- Generally depends on statistical test undertaken
- Can be computed from first principles
- Used to be reflected in statistical tables
  - Compute value
  - Look up p value
- Now Excel, SPSS, etc.

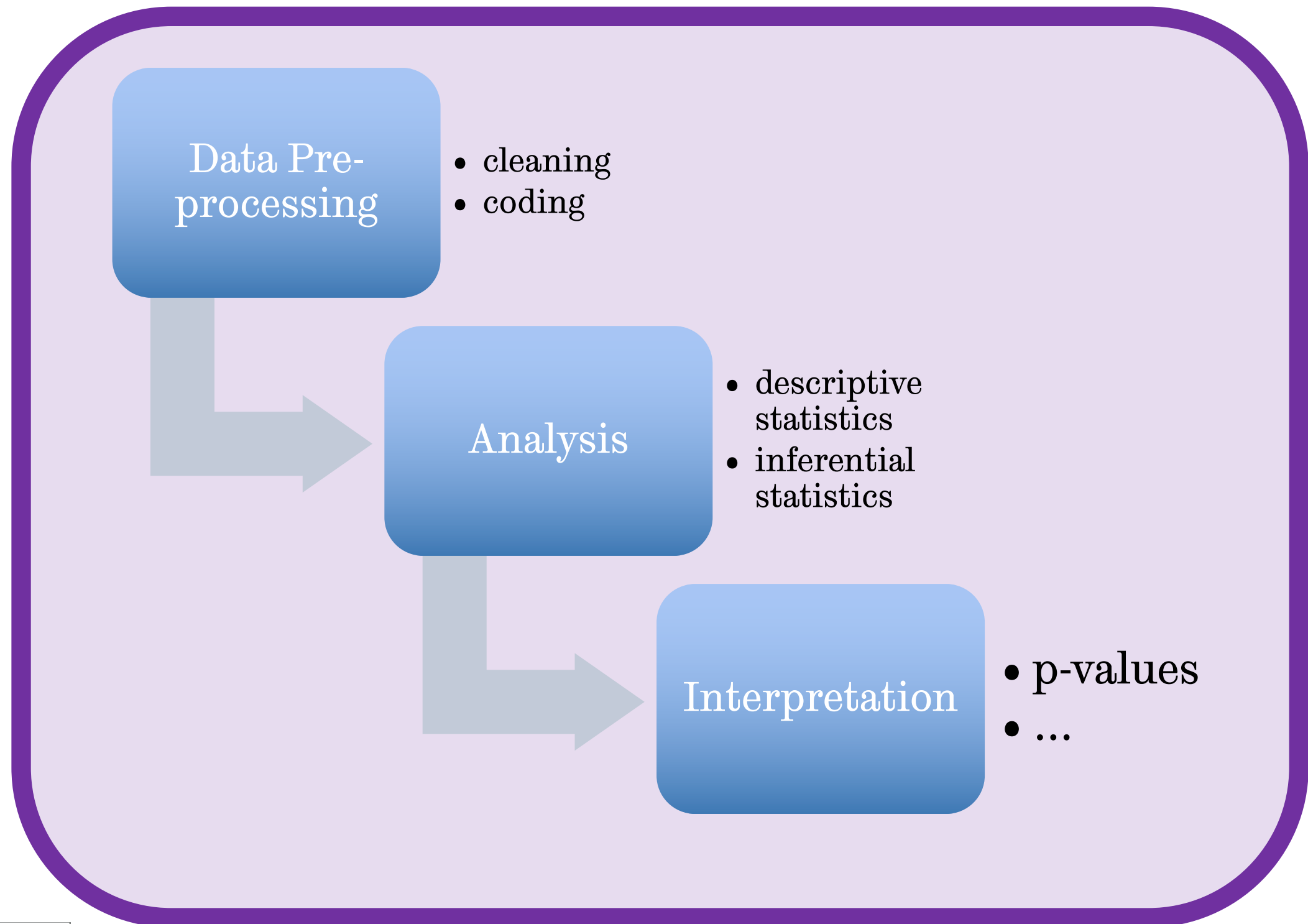


An experiment DOES NOT \*PROVE\* THE HYPOTHESIS, it just provides evidence (statistically significant?) to support the hypothesis.

# Experimental Research Methods

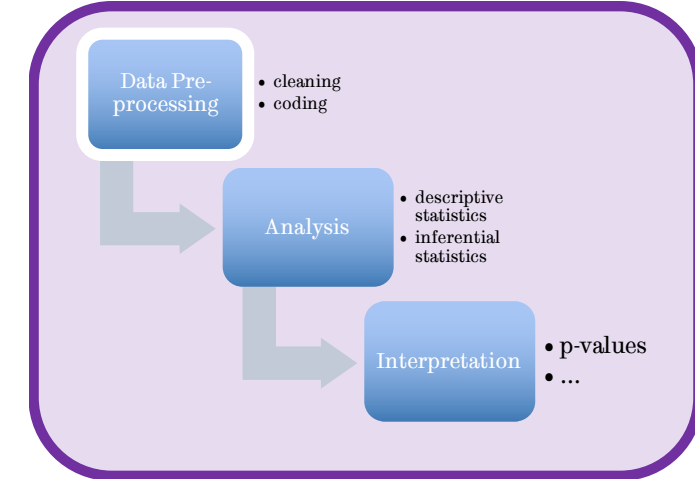


# Data Collection and Analysis



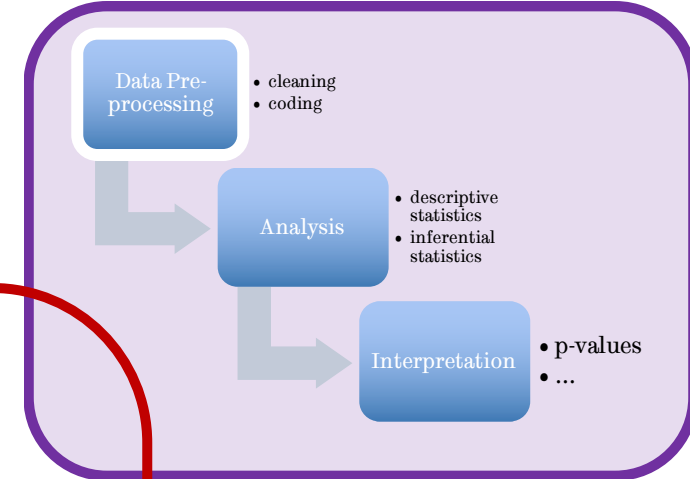
# Data pre-processing

- Cleaning up data
  - Detect errors
  - Formatting
- Coding
  - Types of data that need to be coded
  - Be consistent
- Organizing the data
  - Accommodate to the requirements of statistical software



# Example

Variable	Coding scheme
Demographic characteristics	
Gender: Female	Dichotomous: 0 = no, 1 = yes
Race	
African American	All dichotomous: 0 = no, 1 = yes
Asian American	
Latino	
White	
Age	Range from 16 to 73

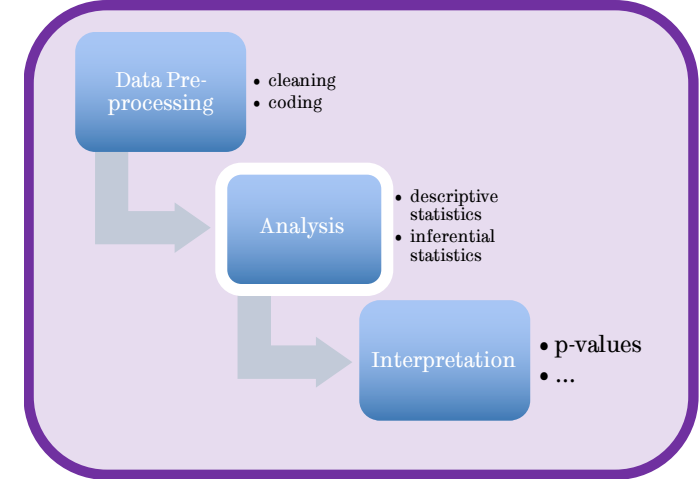


- What happens to missing data?
- What happens with cases where a wrong value is provided? (e.g. age = 270...)

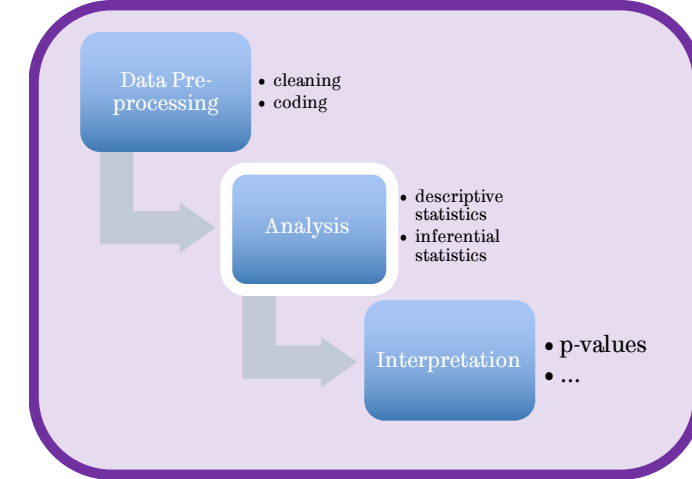
# Descriptive statistics

- Purpose:

- understanding the nature of the data set
- Typically used:
  - Measures of spread:
    - variances, standard deviations, ranges...
  - Measures of central tendency:
    - means, medians, modes...



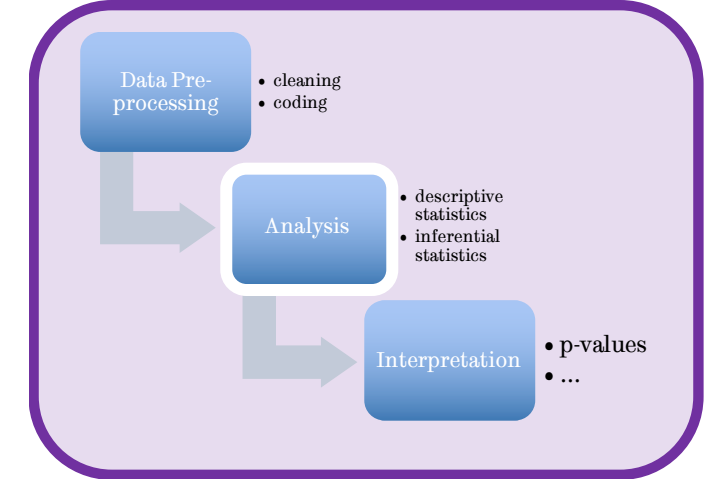
# Example



Variable	Coding scheme	Mean	Standard deviation
Demographic characteristics			
Gender: Female	Dichotomous: 0 = no, 1 = yes	.54	.49
Race			
African American	All dichotomous: 0 = no, 1 = yes	.03	.17
Asian American		.38	.49
Latino		.14	.35
White		.35	.48
Age	Range from 16 to 73	20.29	2.96

# Descriptive statistics

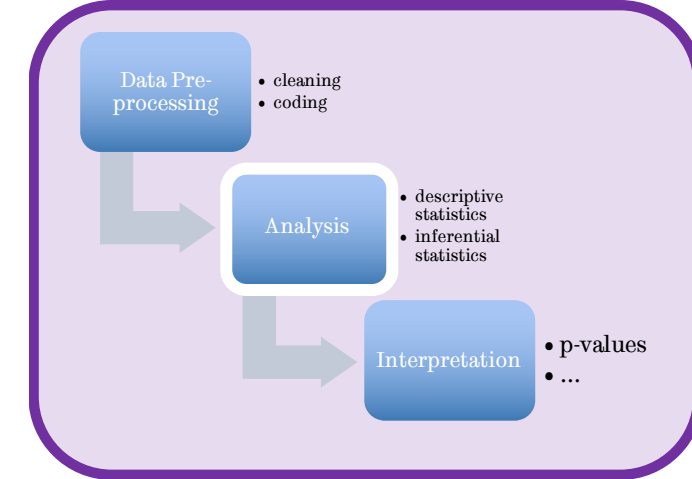
- Measures of spread indicate:
  - How much the data points deviate from the centre of the data set
  - How spread out the data is
  - E.g.:
    - Range, variance and standard deviation
      - Range measures the distance between the highest and lowest scores in the data set



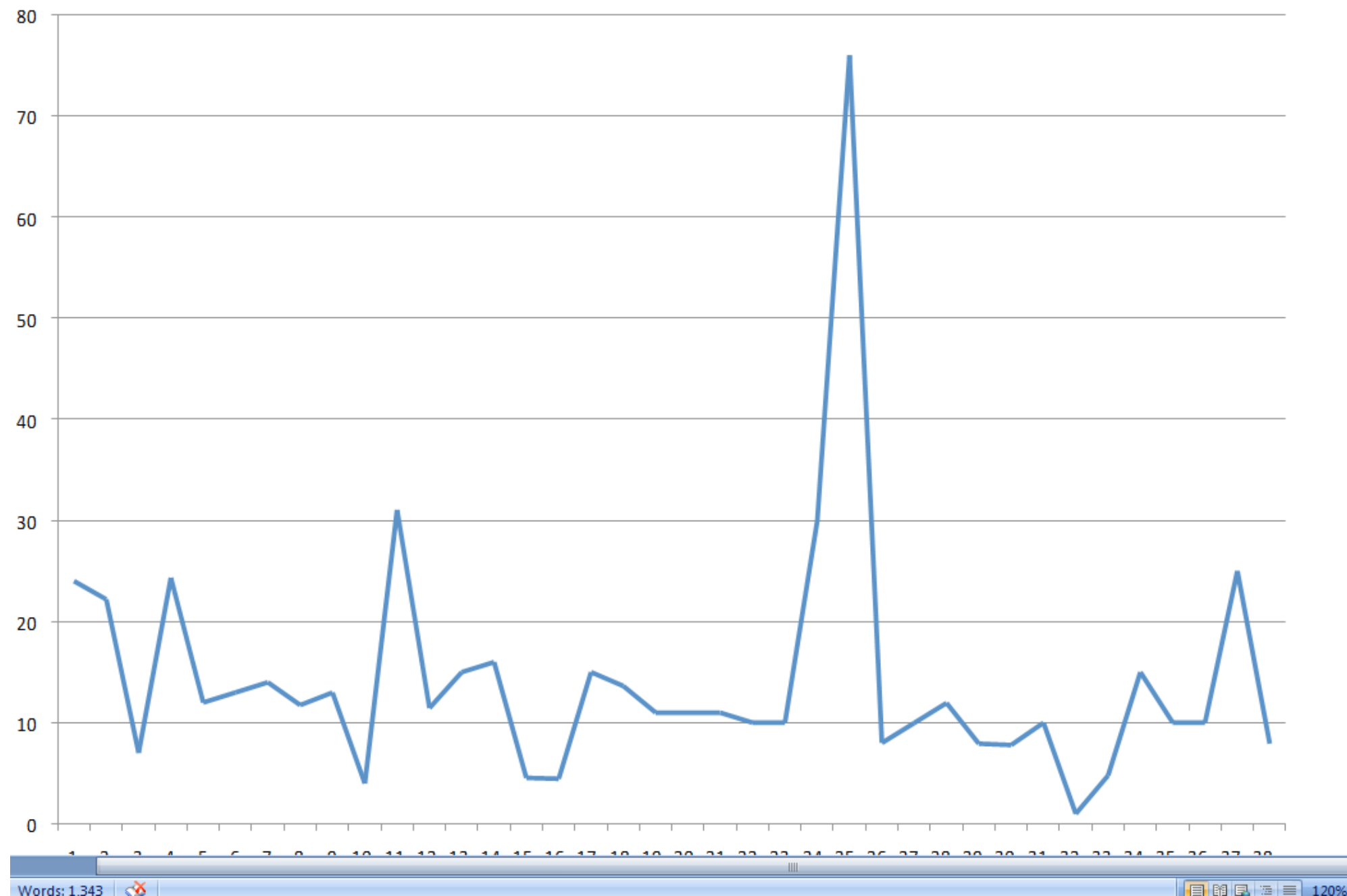
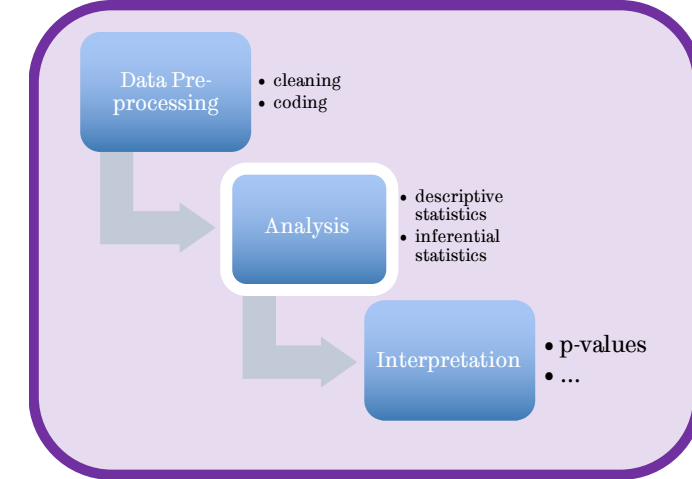


# Descriptive statistics

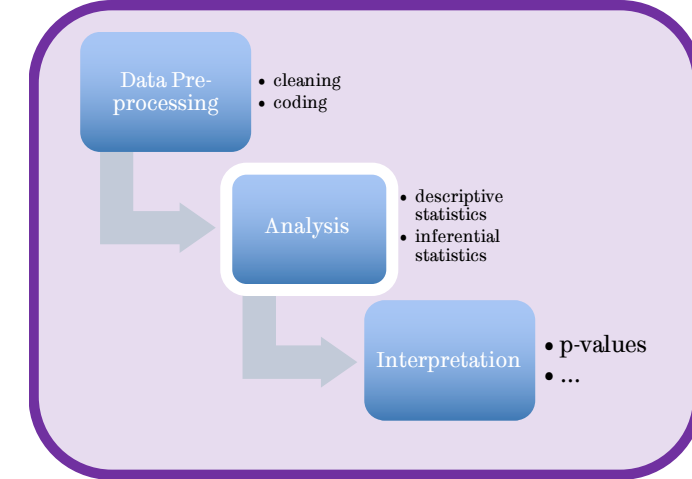
- Measures of central tendency indicate:
  - Where the bulk of the data is located
  - E.g.:
    - Mean of a data set or arithmetic average
  - Can be useful to compare the means of multiple groups
    - If the means are different, could conduct a significance test to see if that difference is statistically significant - see later



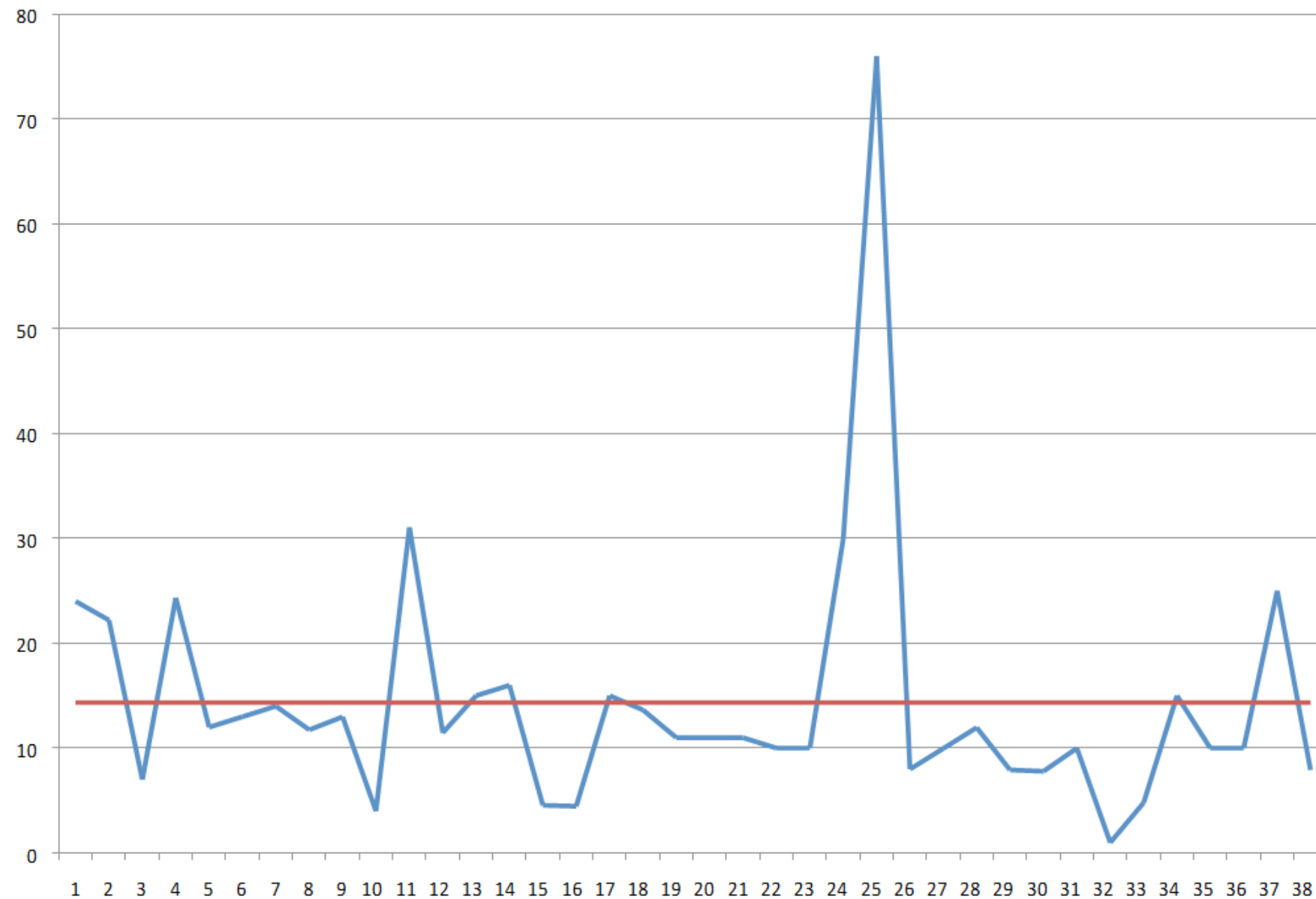
# Descriptive statistics



# Descriptive statistics

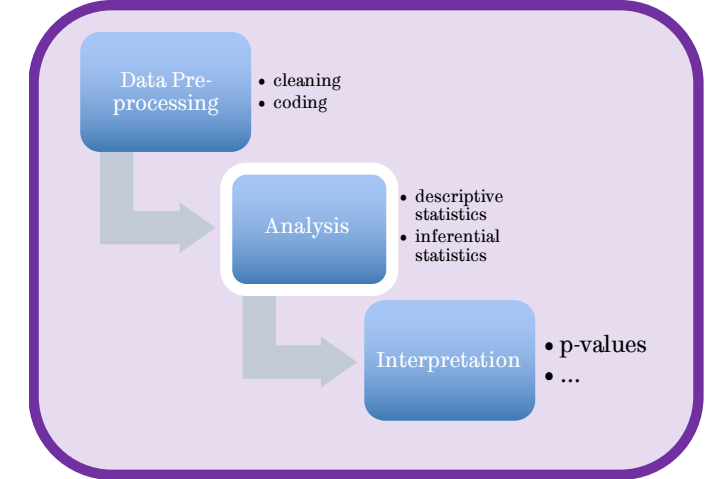


## Mean

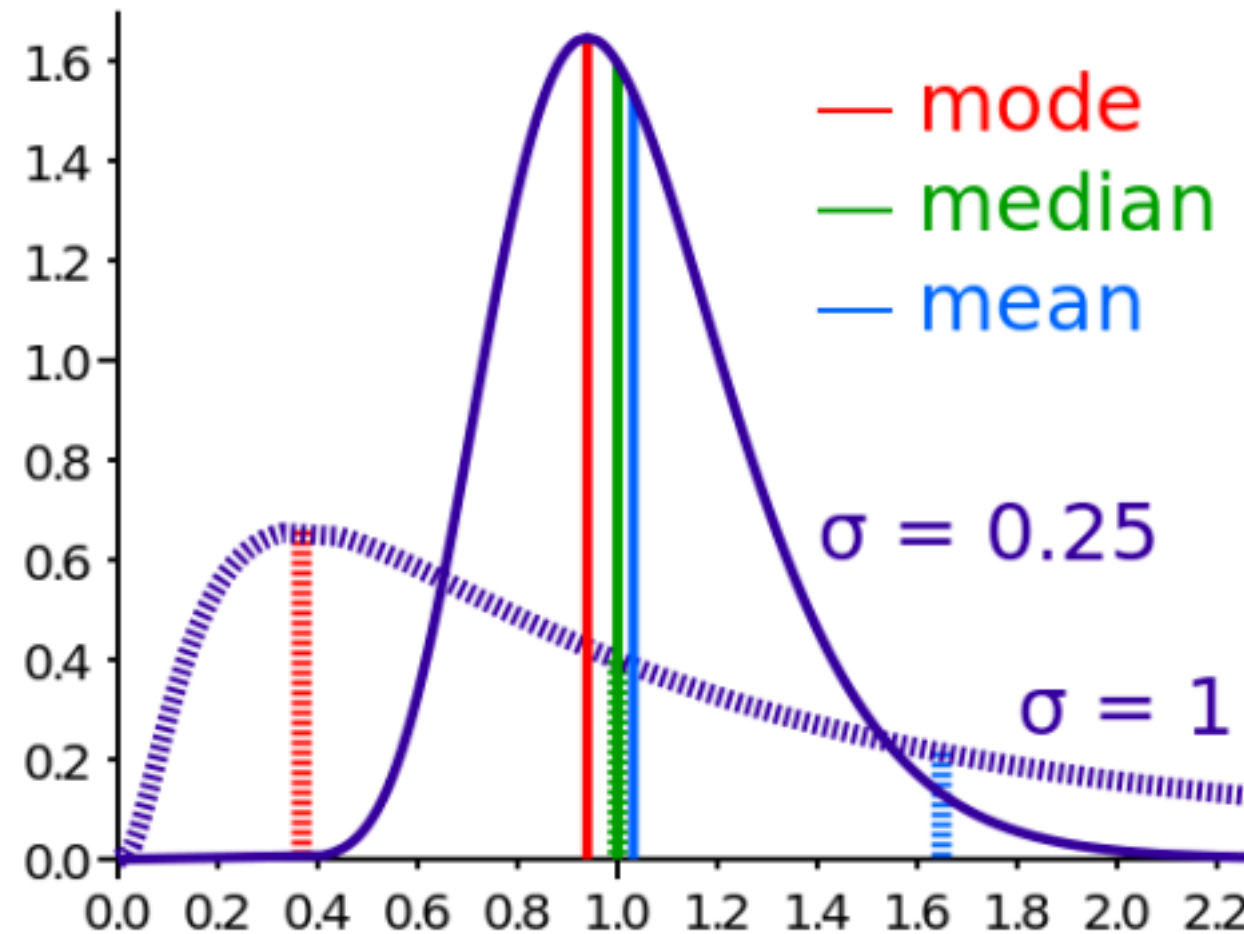
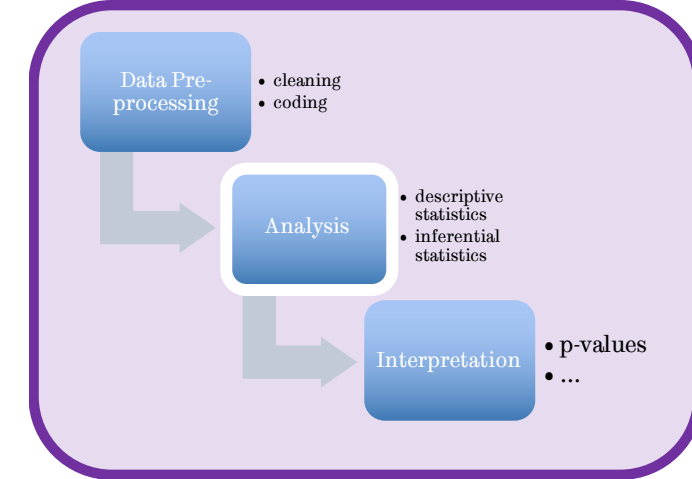


# Descriptive statistics

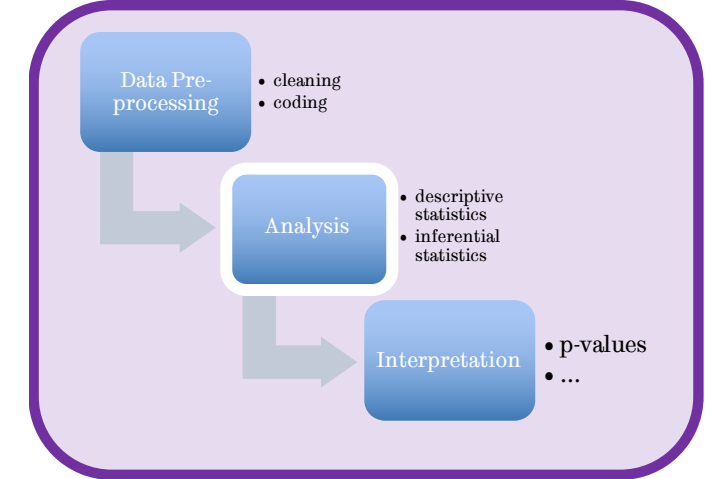
- Other measures of central tendency
- Median
  - Order values, median is central value (or mean of two central values)
  - If data is skewed then mean and median will differ significantly
- Mode
  - The most commonly occurring value
    - Rarely used



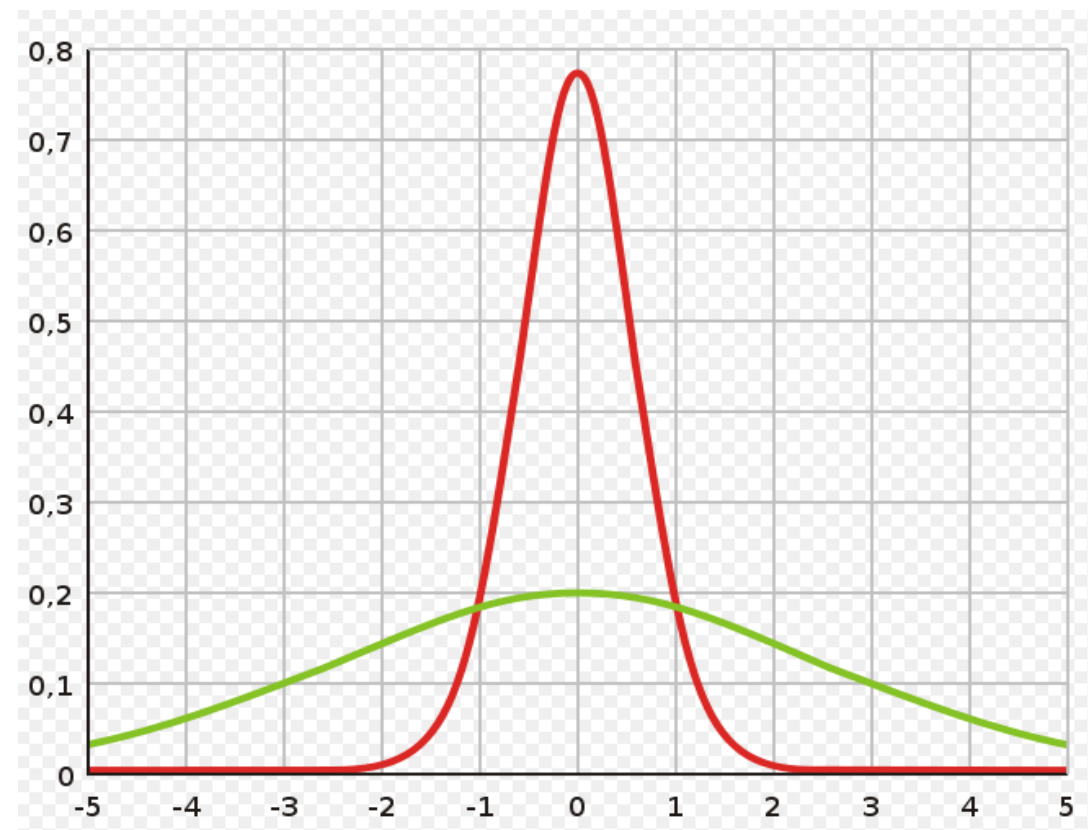
# Descriptive statistics



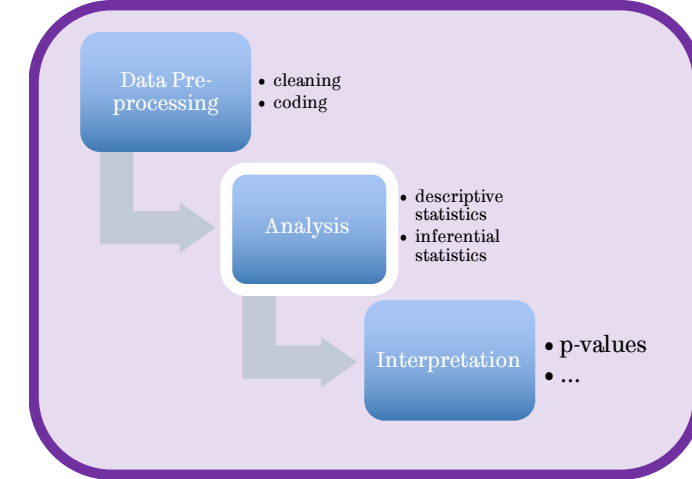
# Descriptive statistics



- The variance of a data set is the mean of the squared distances of all the scores from the mean of the data set
- The square root of the variance is called the standard deviation

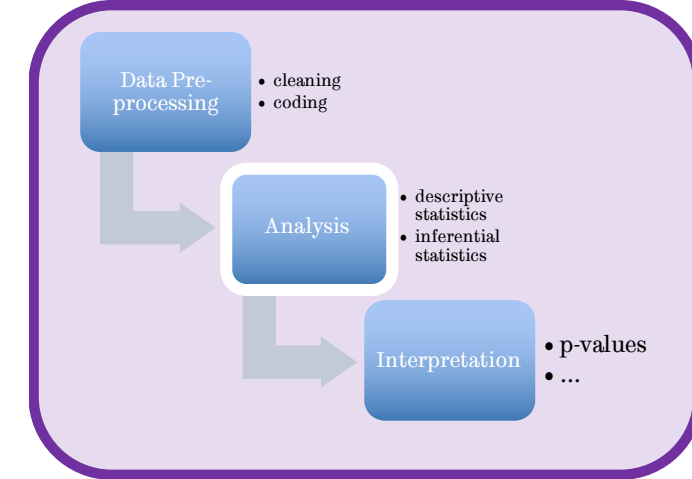


# Descriptive statistics

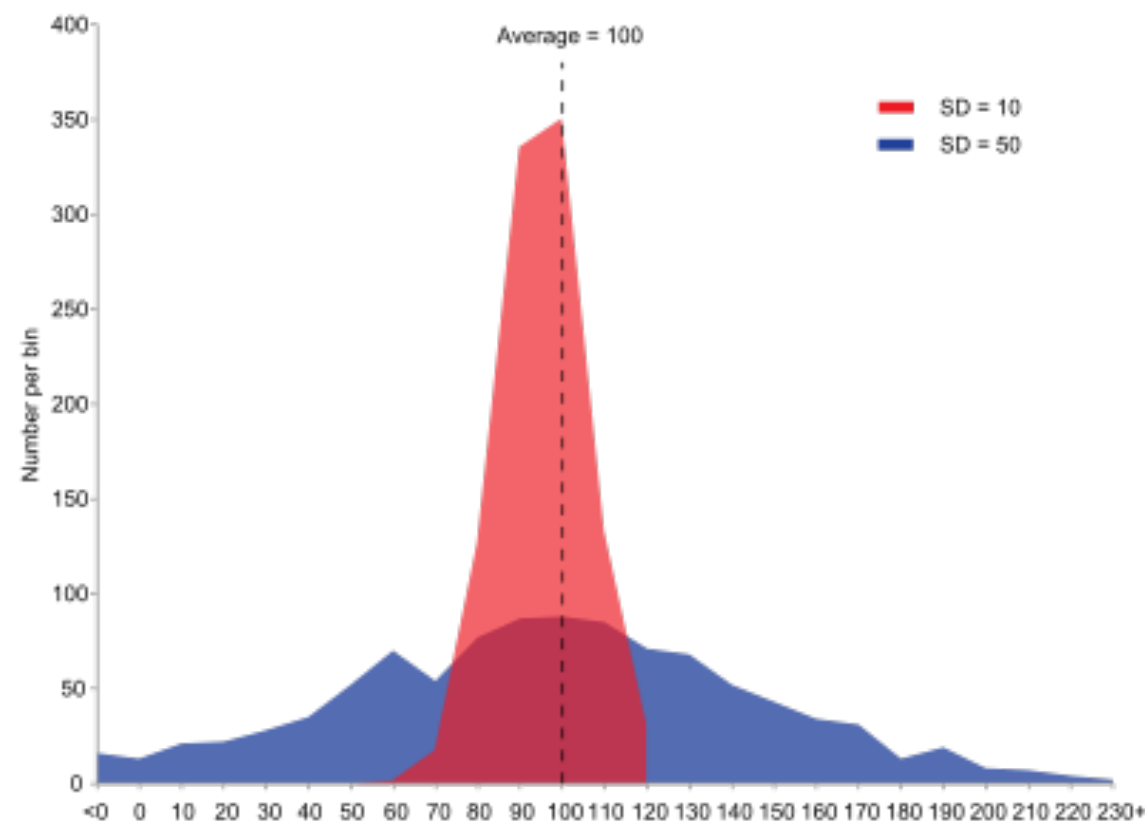


- Variance:
  - measure of how far a set of numbers is spread out
  - zero-variance indicates that all the values are identical.
  - non-zero variance is always positive:
    - a small variance indicates that the data points tend to be very close to the mean (expected value) and hence to each other,
    - while a high variance (or standard deviation) indicates that the data points are very spread out from the mean and from each other.

# Descriptive statistics

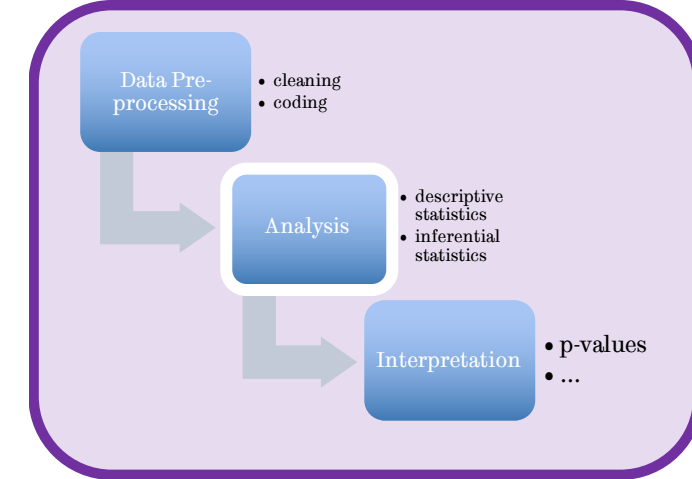


- Two sample populations with the same mean and different standard deviations

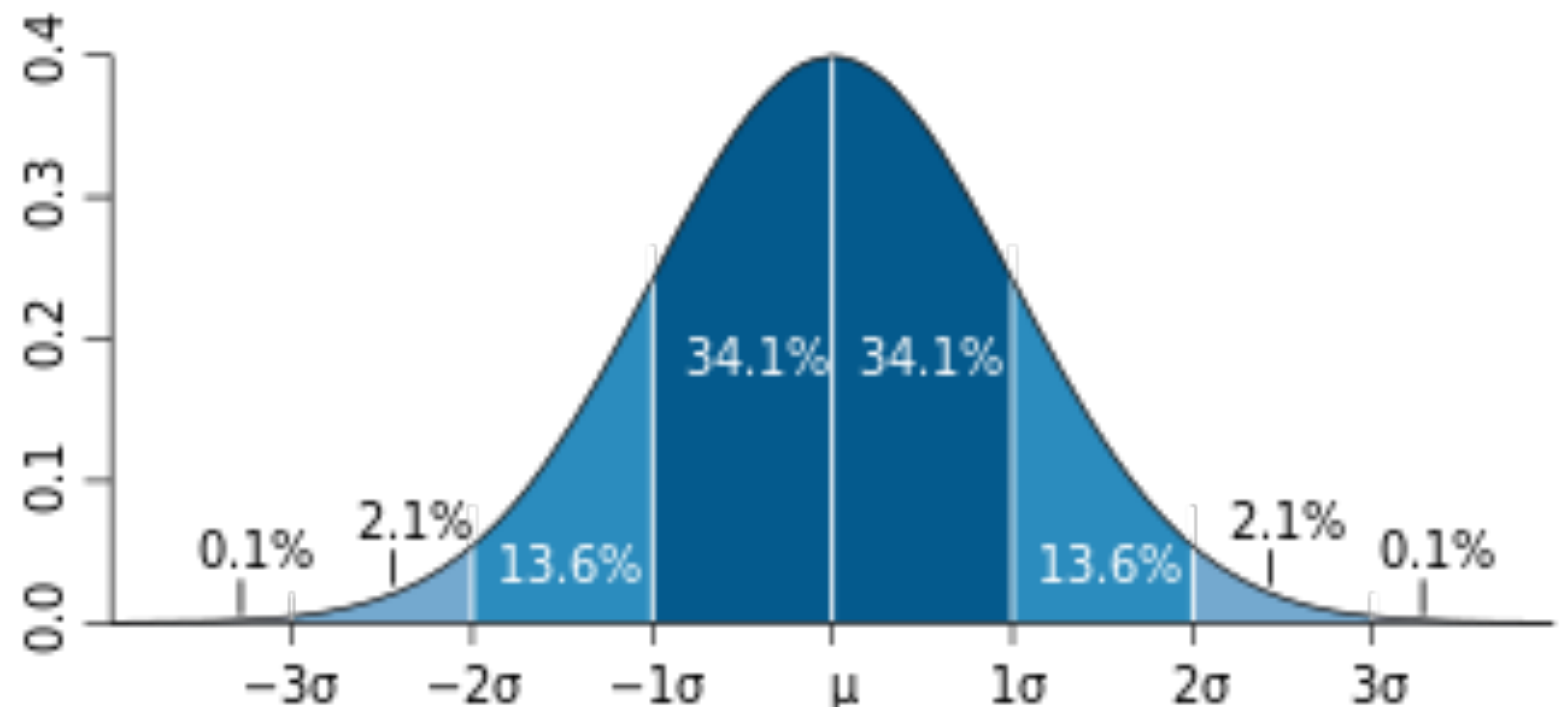




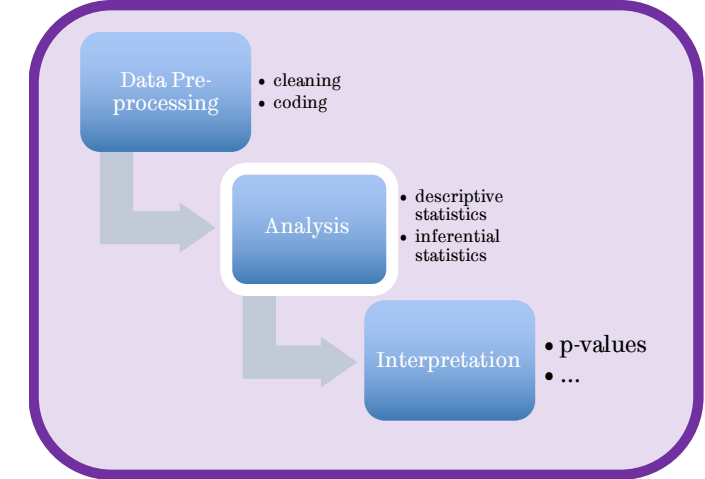
# Descriptive statistics



- Plot of a normal distribution (or bell-shaped curve) where each band has a width of 1 standard deviation
  - many attributes from different fields are normally distributed
    - e.g. heights of a population, student grades, IQ
  - Parametric tests assume that the data is normally distributed



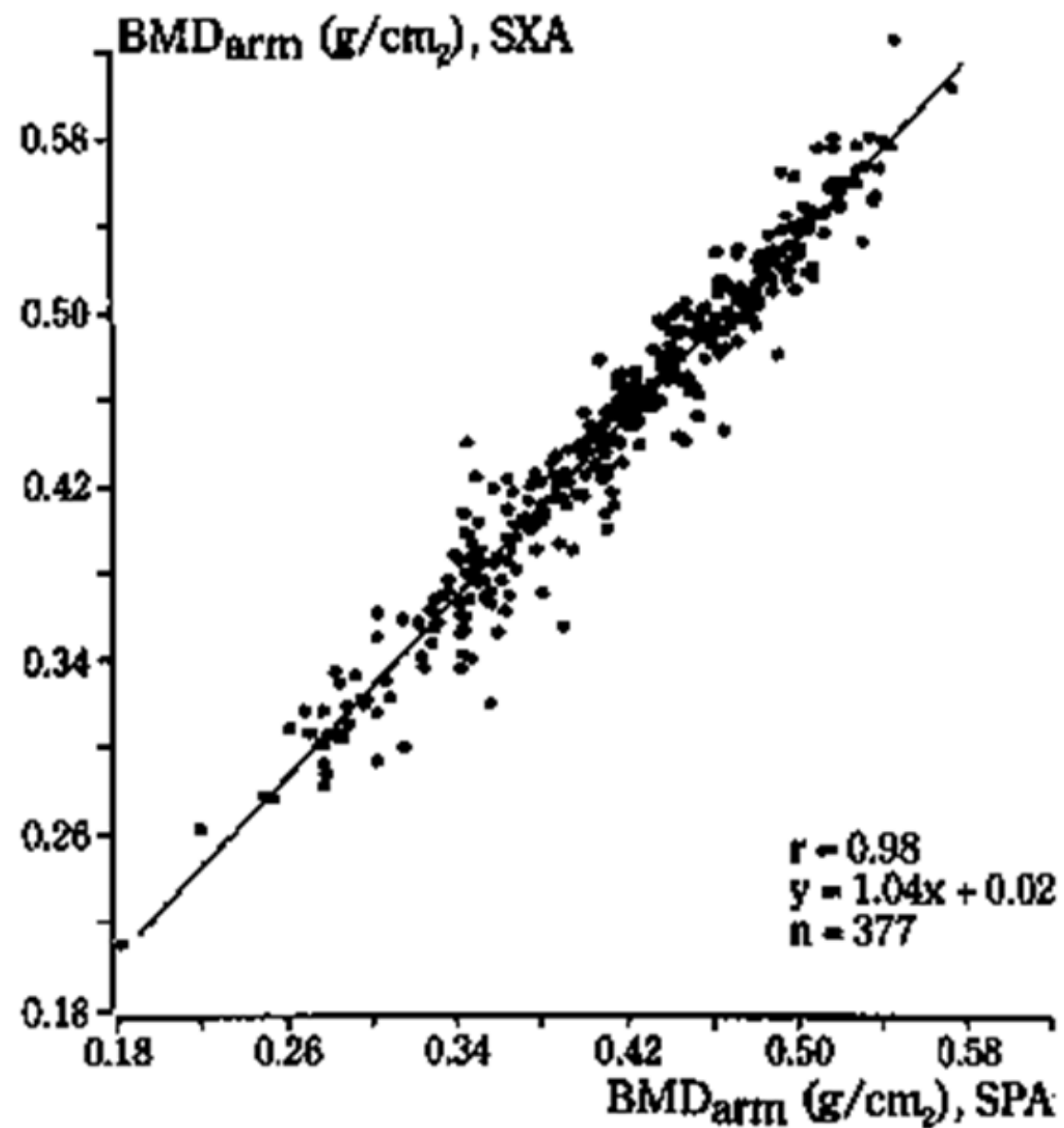
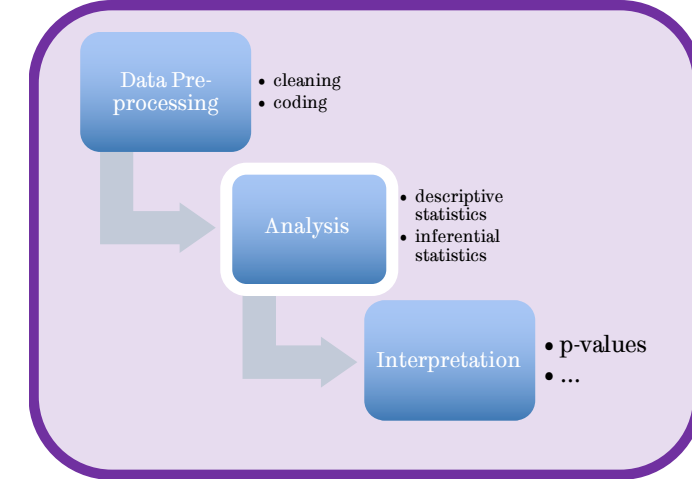
# Descriptive statistics



- Correlation
  - Measure of how changes in variable X match changes in variable Y
  - E.g.
    - Distance to work correlates with time taken to reach work
    - Level of education correlates with salary

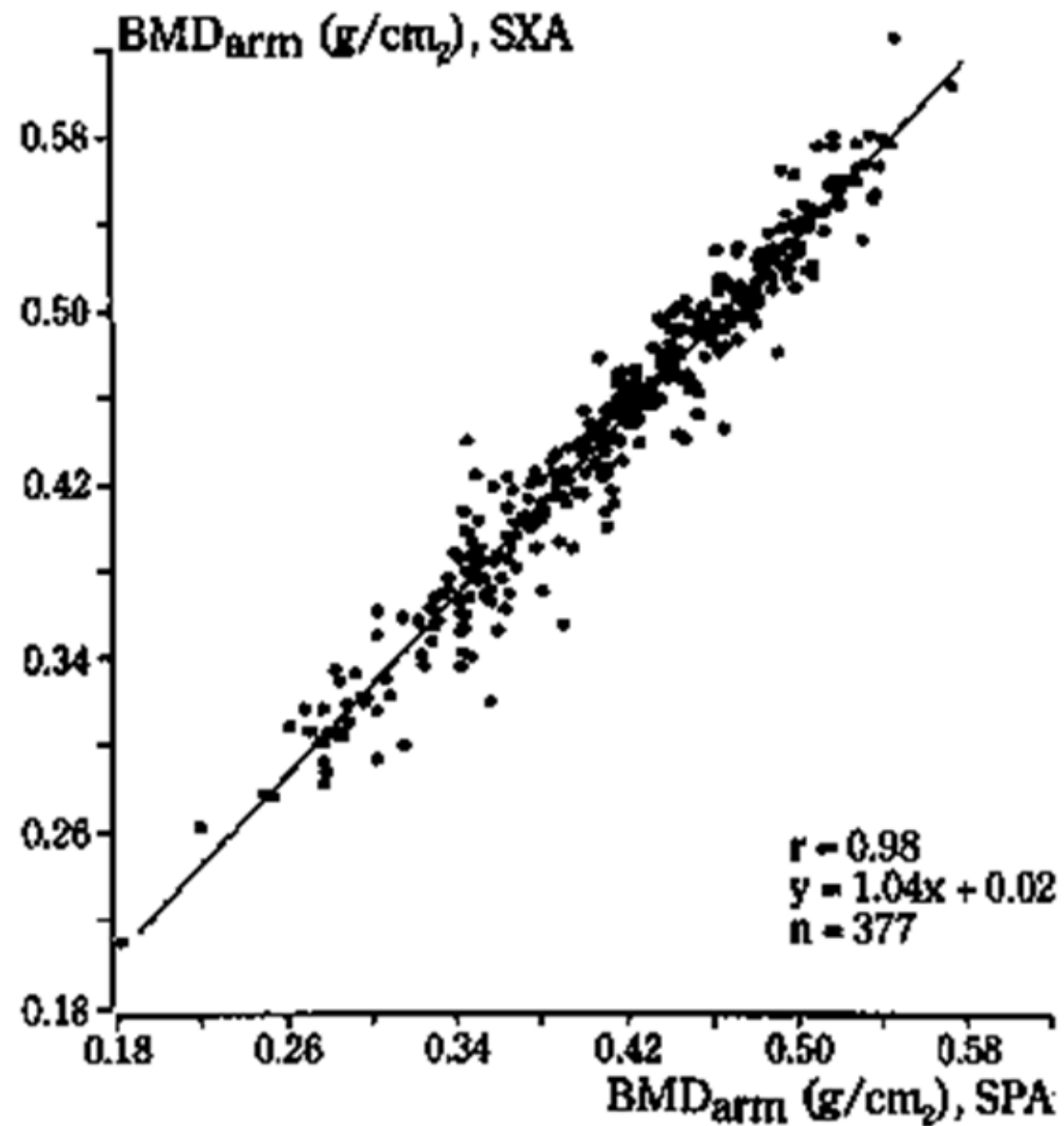
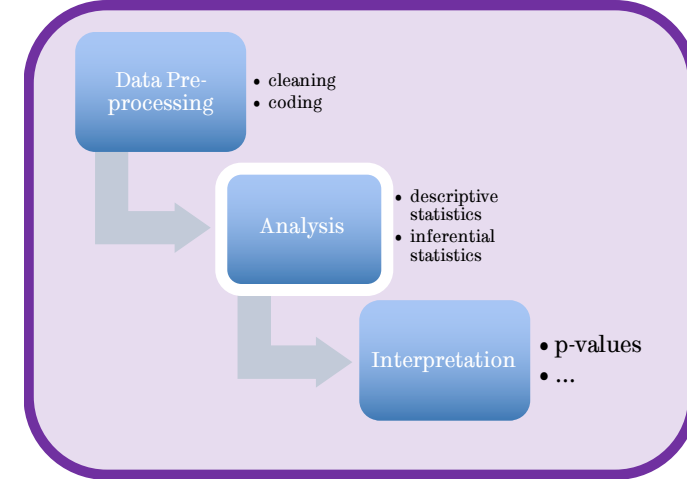
# Task

- Correlated or uncorrelated?



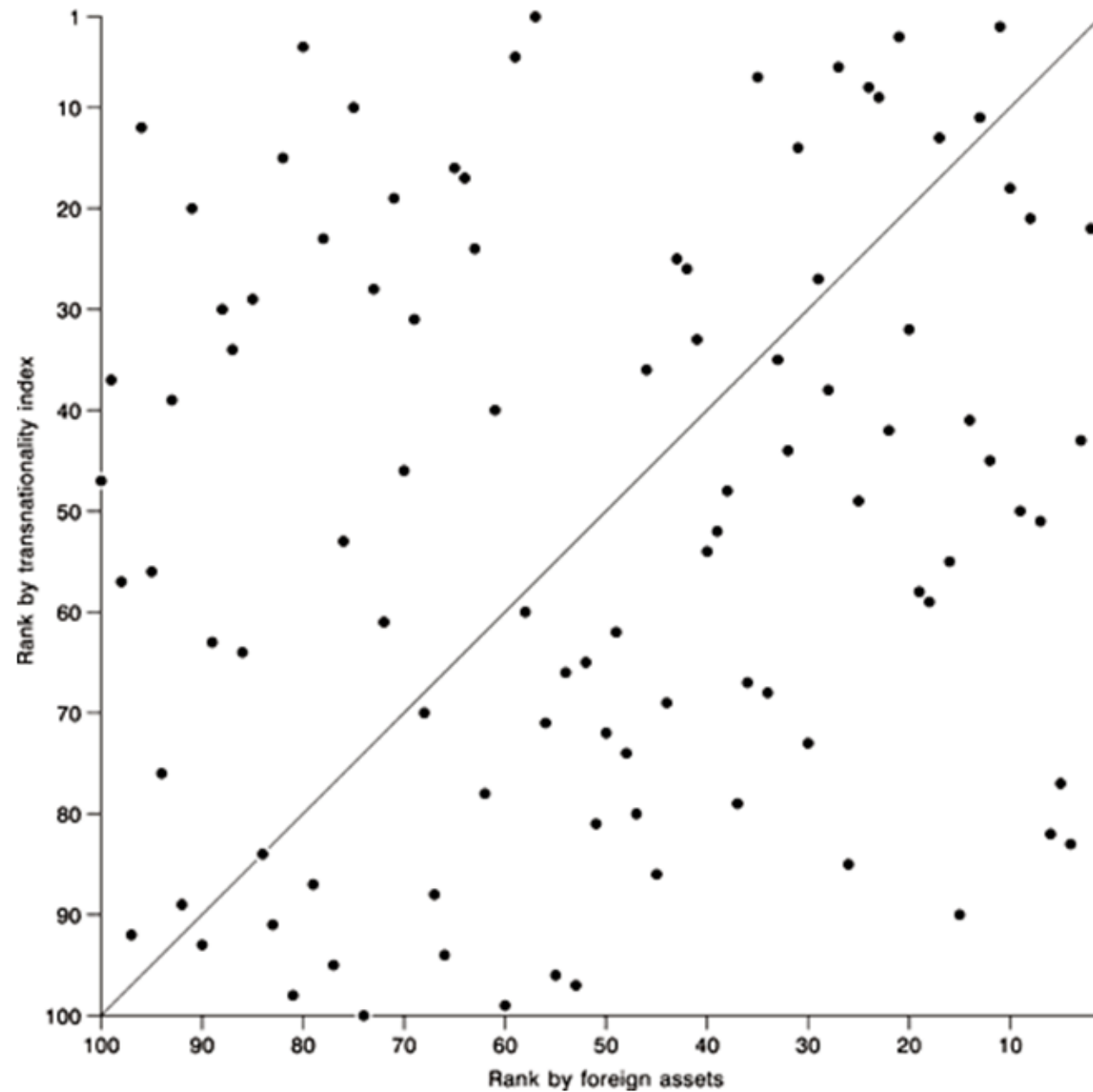
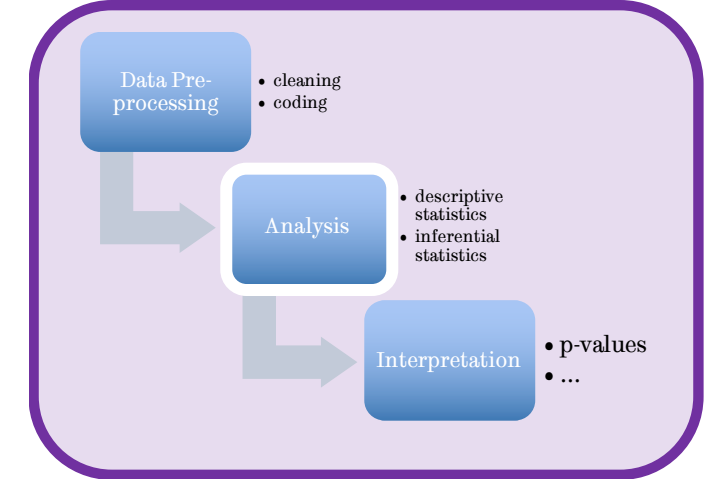
# Example

- Highly correlated



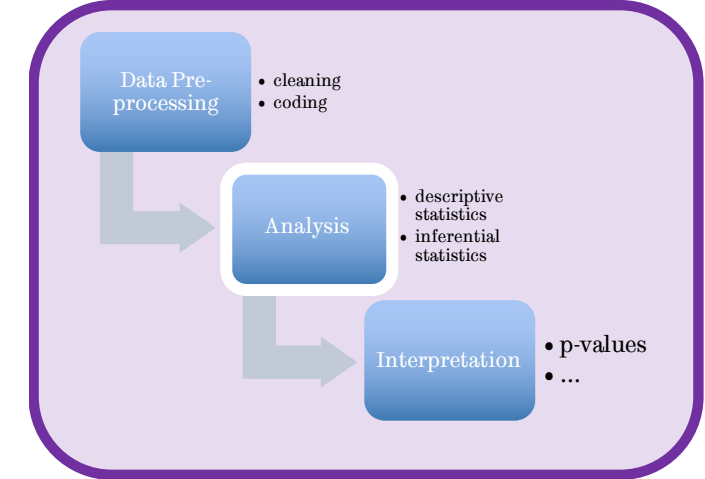
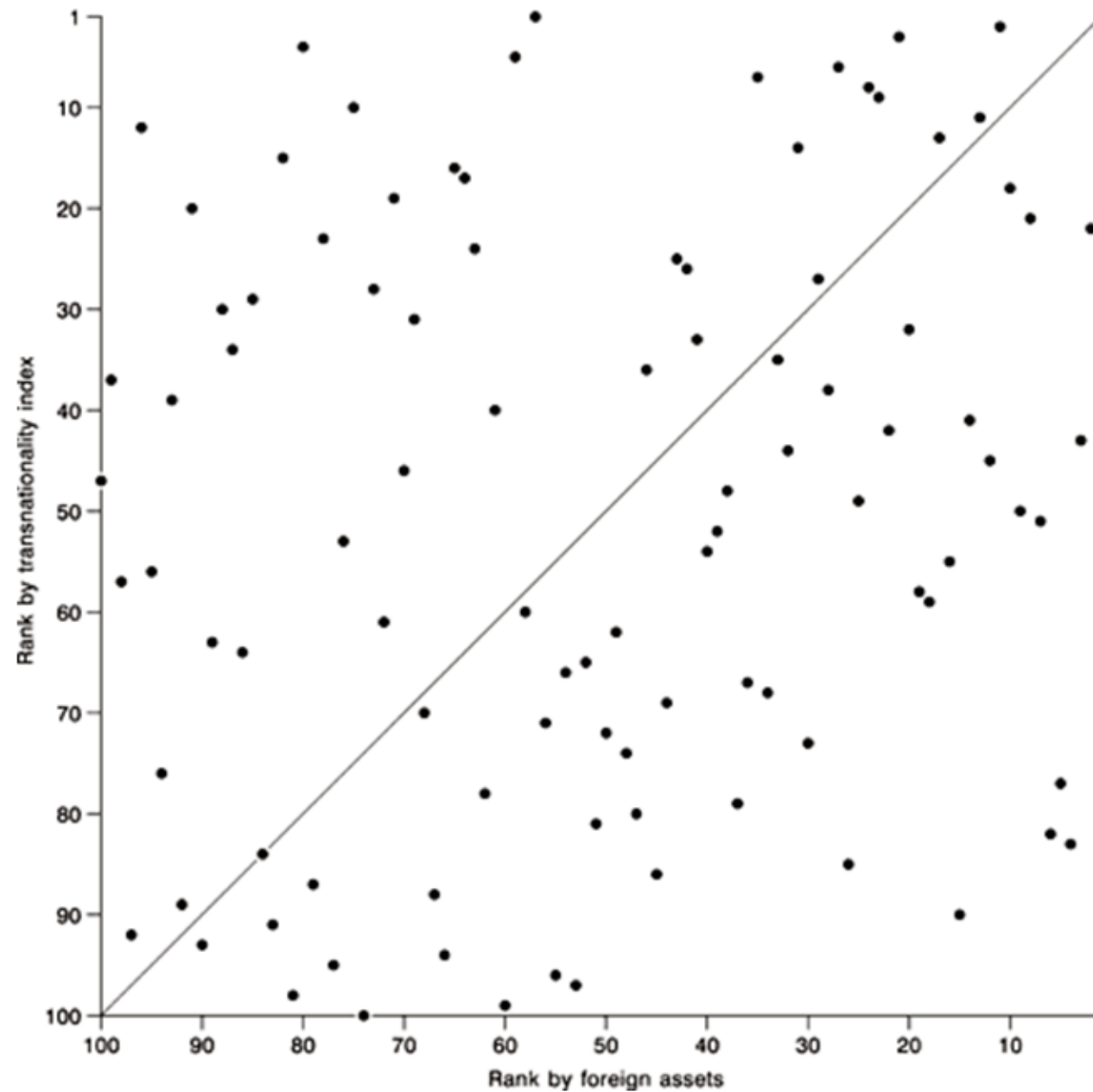
# Task

- Correlated or uncorrelated?



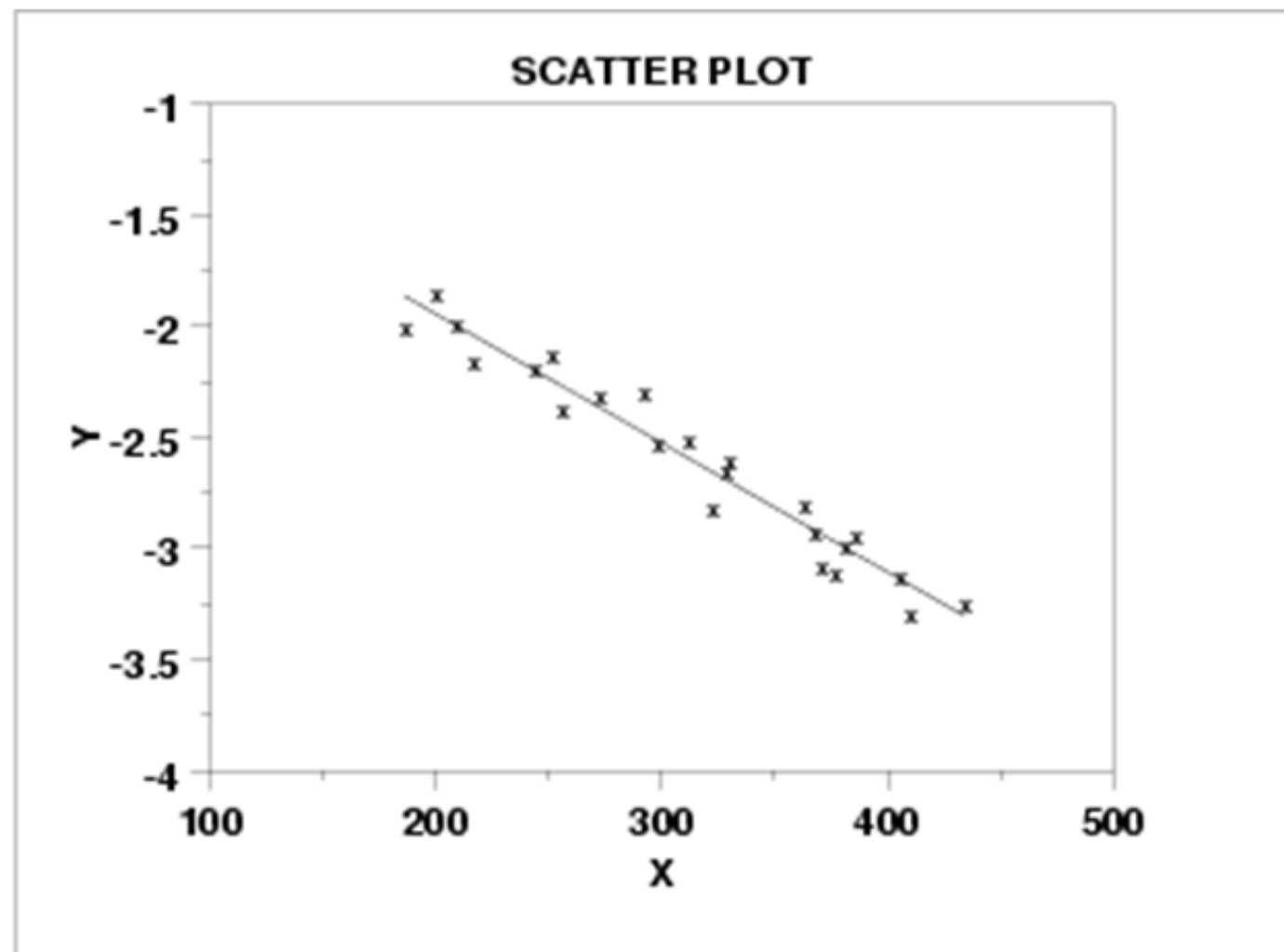
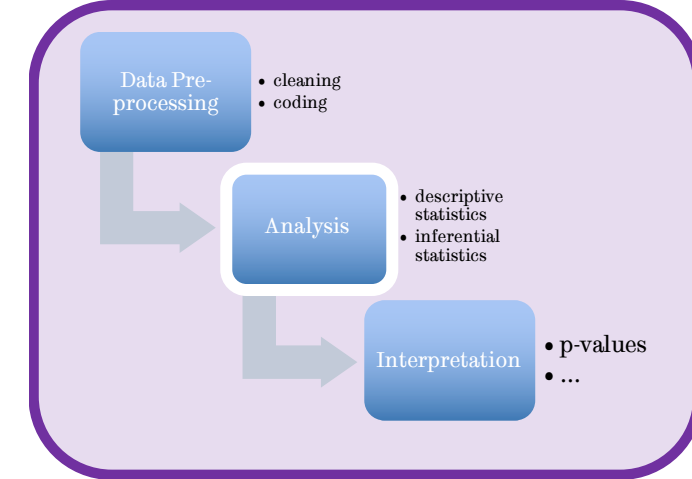
# Example

- Uncorrelated



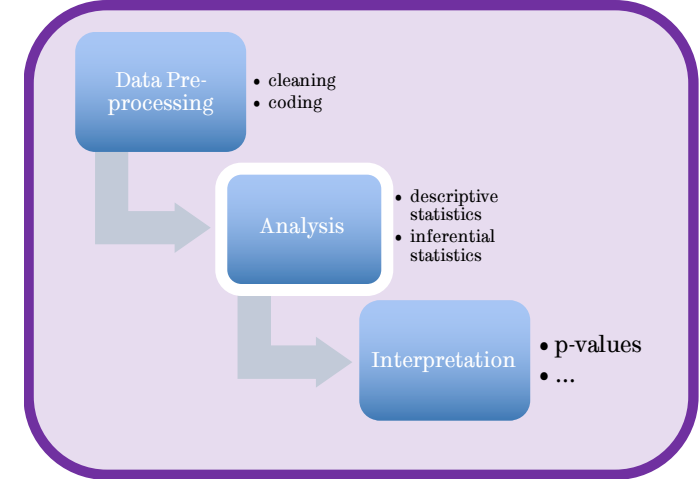
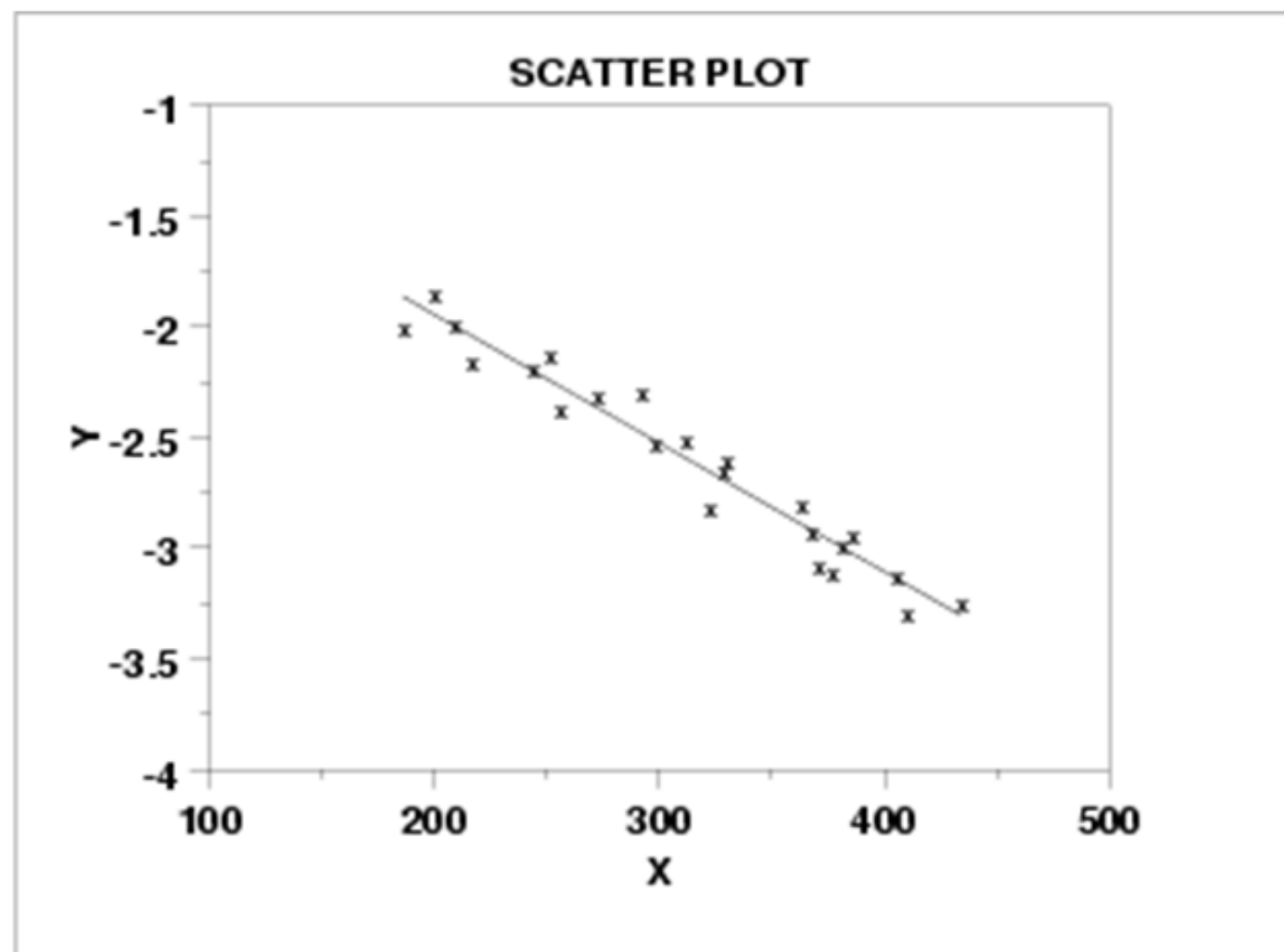
# Task

- Correlated or uncorrelated?



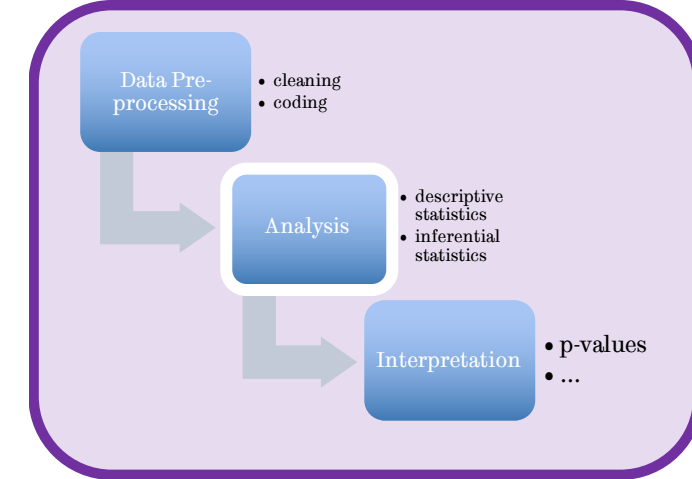
# Example

- Negatively correlated



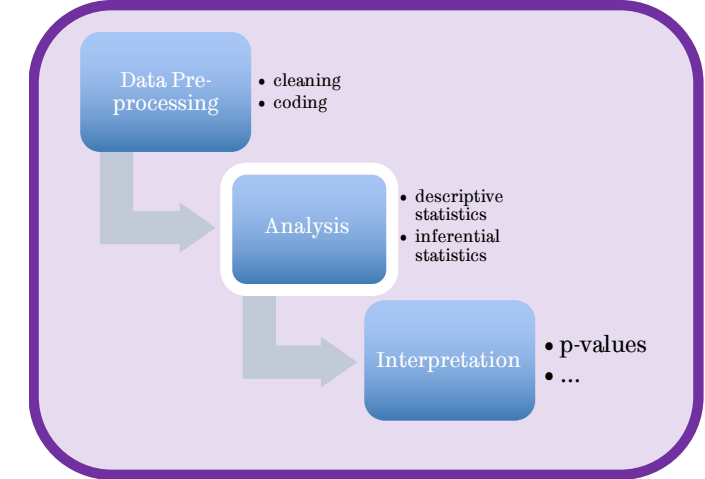


# Descriptive statistics



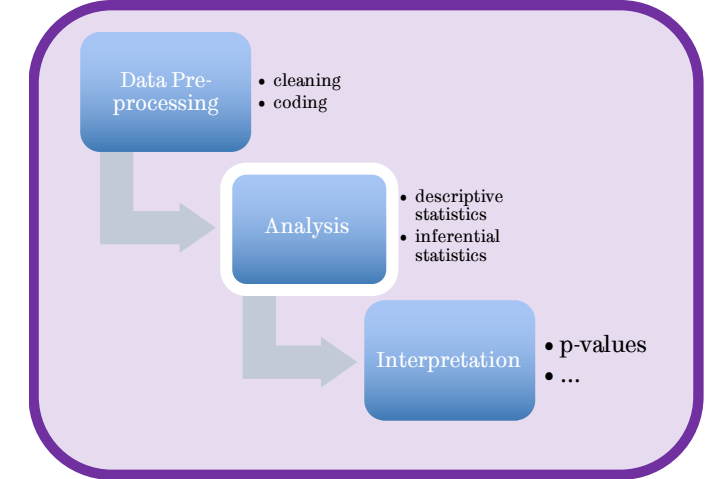
- Correlation:
  - tells us very little
  - suggests a relationship between two variables
  - relationship which might not necessarily be linear
  - does not suggest that there is a direct connection
- Correlation is NOT causation
  - Classic scientific fallacy
  - Because two variables move in the same direction does not mean they are linked
    - Secondary factor
    - Weird correlation effects

# Task



- Correlation between women taking Hormone Replacement Therapy (HRT) and decreased Coronary Heart Disease.
- But Controlled Trials demonstrated that HRT significantly increases chances of Coronary Heart Disease
- Why?

# Task

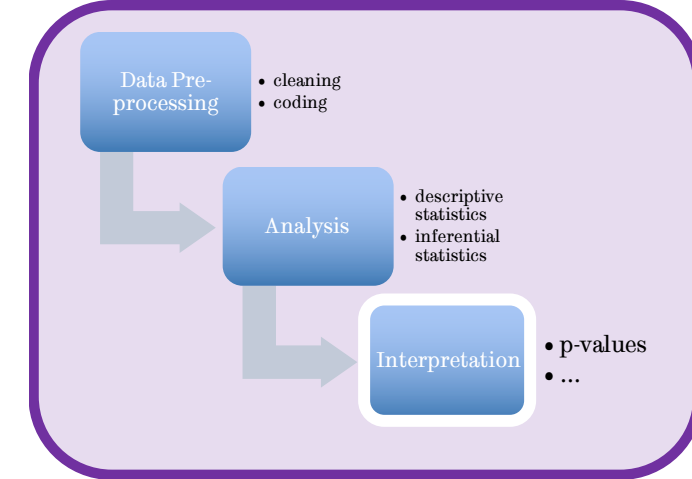


- Correlation between women taking Hormone Replacement Therapy (HRT) and decreased Coronary Heart Disease.
- But Controlled Trials demonstrated that HRT significantly increases chances of Coronary Heart Disease
- Why?

**confounding factor**

# Descriptive statistics

- Examining the data
  - Can reveal nature of relationship
  - Can indicate if significant correlation is really significant
- Consider the nature of effects
  - Why is there correlation
  - Is there a confounding factor?
  - Common sense

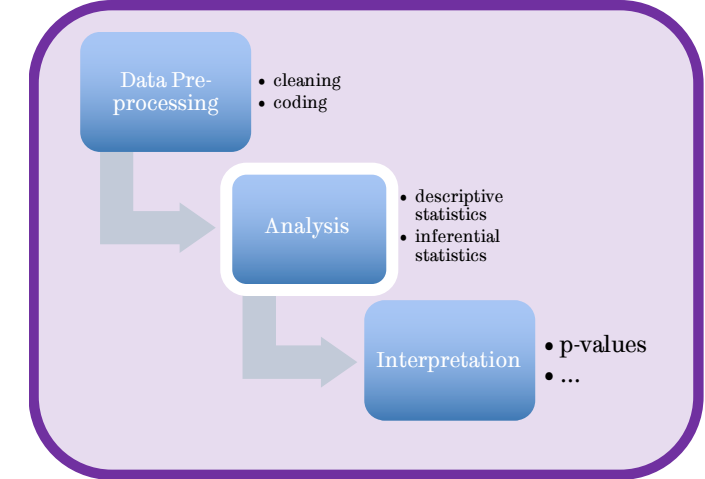


# Experimental research

- We measure:
  - Dependent Variables
- In the presence of
  - Independent Variables

# Inferential Statistics

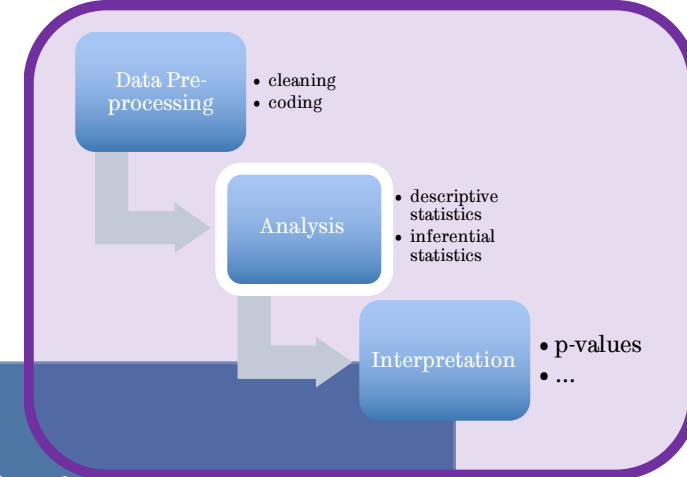
- Tests of significance



# Statistics books

- Usually the thinner the better
- 
- Bigger Books
  - Go through everything from first principles
- Thinner Books
  - Stick to the things you'll use professionally

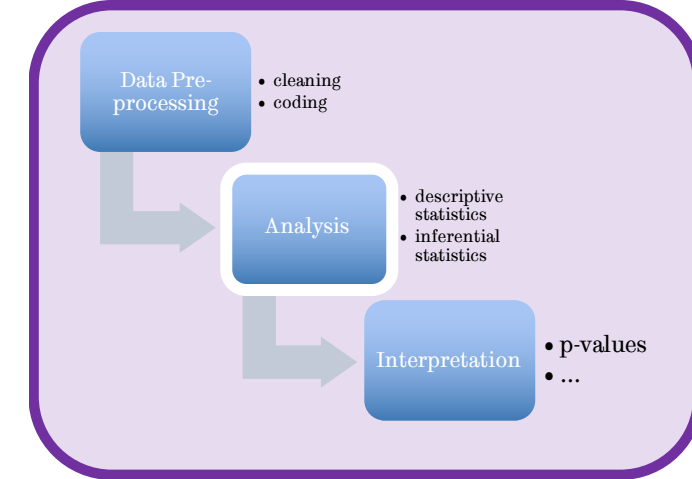
# Comparison of means



Experiment design	Number of Independent variables	Number of conditions for each independent variable	Types of significance tests
Between group	1	2	Independent samples t test
	1	3 or more	One way ANOVA
	2 or more	2 or more	Factorial ANOVA
Within group	1	2	Paired samples t test
	1	3 or more	Repeated measures ANOVA
	2 or more	2 or more	Repeated measures ANOVA



# Significance tests



- Choosing the right statistical test for the experiment is not straightforward:
  - there are more tests available depending on the nature of the samples (e.g. non parametric tests such as Mann-Whitney U test for differences between independent samples, and Sign Test for differences between related samples)
- Experiment design is far more important than statistical accuracy
  - T-tests for comparing means in two conditions
  - ANOVAs for more than two conditions
  - Correlation for comparing two populations
  - CHI Squared when conditions are right

# Summary

- Quantitative research methods
  - Null hypothesis and alternative (or experimental) hypothesis
  - Significance, p-values, confidence interval...
  - Descriptive statistics vs inferential statistics
    - Measures of central tendency
    - Measures of spread
    - Correlation
  - Tests of significance and their suitability

# Evaluation in HCI

