

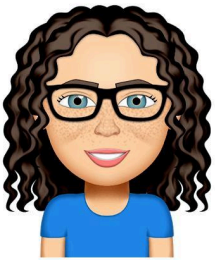
# Level 3

# Inference

# Workbook

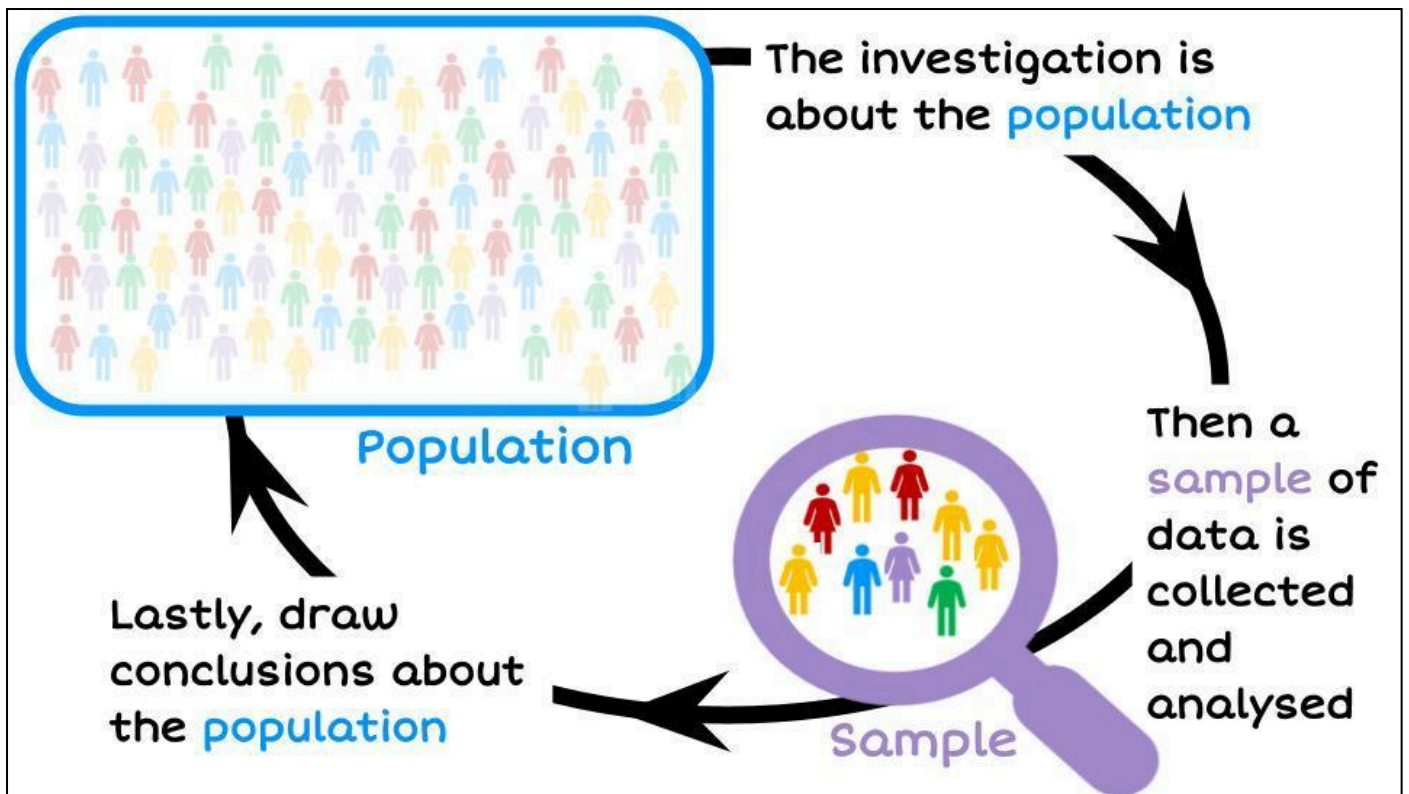


**Name:**

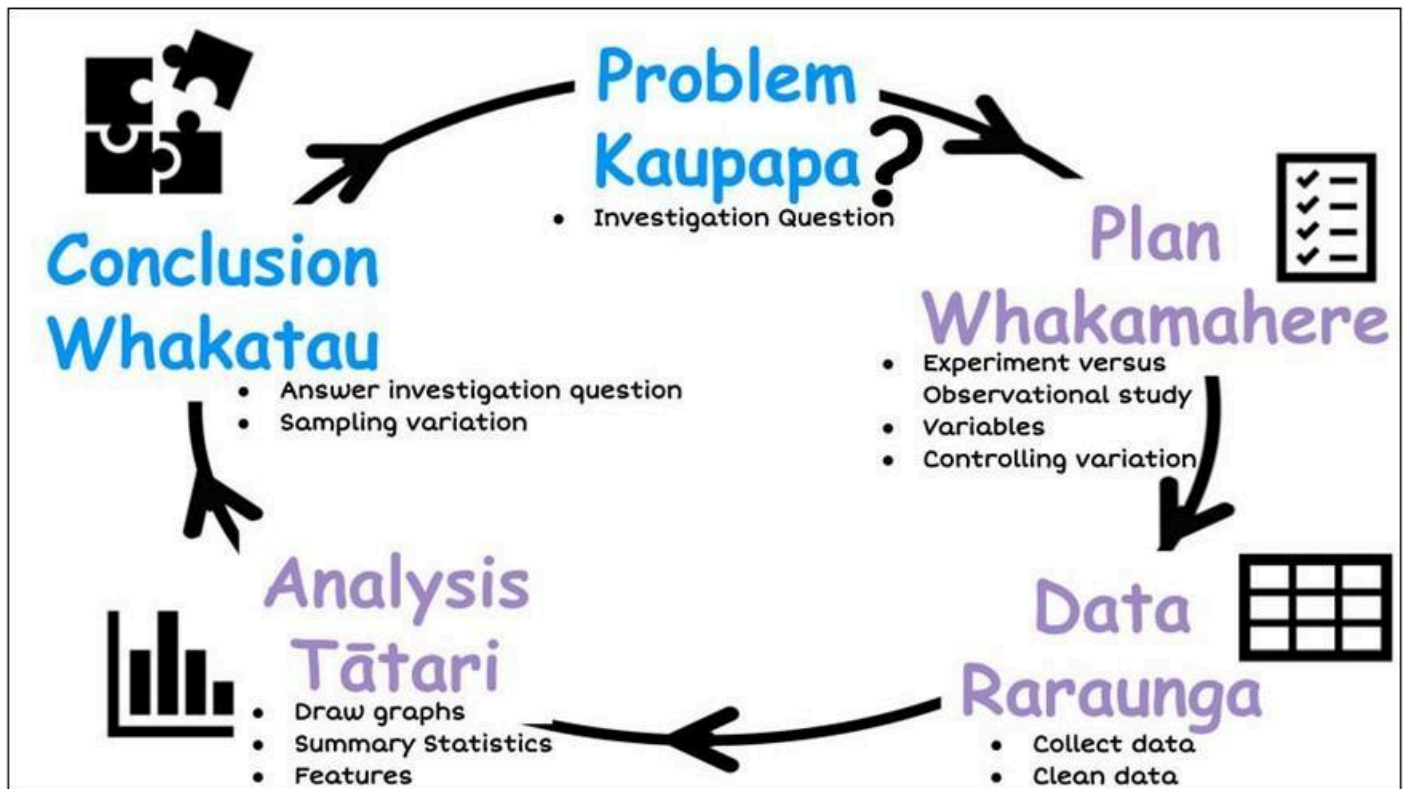


By Liz Sneddon

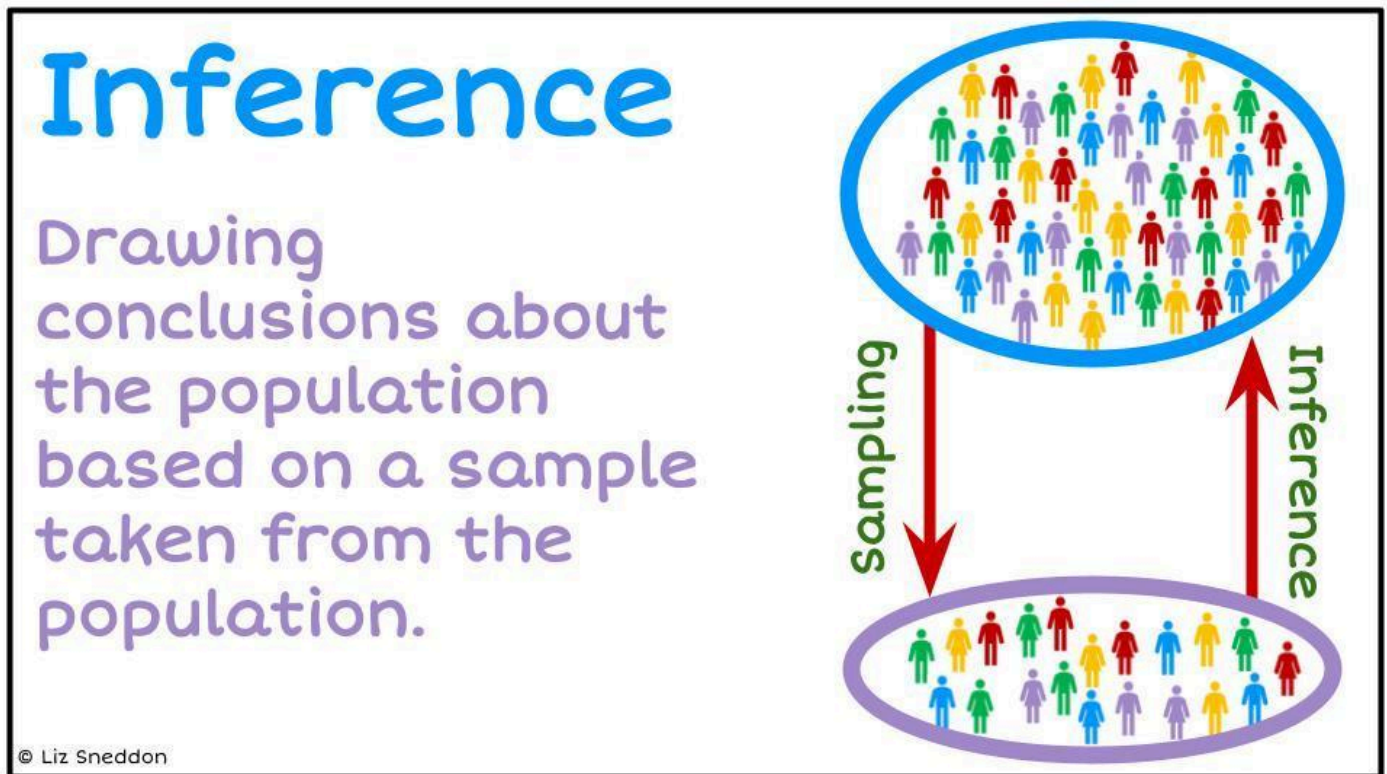
# Investigative Process



The PPDAC cycle is the core of all statistical investigations.



At the end of this investigation, you will need to make an Inference.



## Populations and samples

We start with an investigation question about a population. We often have a hypothesis or prediction of what we expect to find.



# Language

## Language

Unless we have population data (e.g. a census) the results can only ever be **suggestive or inferential**.

Do **NOT** use definitive language (E.g. ~~Prove~~)

© Liz Sneddon



## Average

The NCEA statistics glossary<sup>1</sup> defines the average:

A term used in two different ways.

When used generally, an average is a number that is representative or typical of the centre of a set of numerical values. In this sense, the number used could be the *mean* or the *median*. Sometimes the mode is used. This use of average has the same meaning as *measure of centre*.

When used precisely, the average is the number obtained by adding all values in a set of numerical values and then dividing this total by the number of values. This use of average has the same meaning as *mean*.

As you can see, there are two different definitions of the word average. Therefore, in order to avoid confusion, **please do not use the word average**.

Instead use the specific terminology of **mean, median or mode** as appropriate.

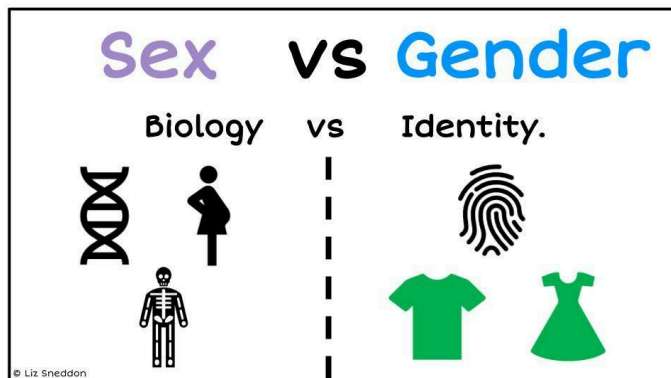
## Gender versus Sex for data collection

<sup>1</sup> <https://seniorsecondary.tki.org.nz/Mathematics-and-statistics/Glossary/Glossary-page-A>

Statistics New Zealand (Tatauranga Aotearoa) is the government organisation that runs the Census (every 5 years). In the 2023 Census, they changed the questionnaire to include questions on both sex and gender, in order to collect data that represents the diversity of Aotearoa New Zealand, along with more accurate and detailed information across population groups.

Here are updated definitions that Statistics NZ has recently released:<sup>2</sup>

**gender:** Refers to a person's social and personal identity as male, female, or another gender such as non-binary. Gender may include the gender that a person internally feels ('gender identity'), and/or the gender a person publicly expresses ('gender expression') in their daily life. A person's current gender may differ from the sex recorded at their birth and may differ from what is indicated on their current legal documents. A person's gender may change over time. Some people may not identify with any gender.



**sex at birth:** Refers to the sex recorded at a person's birth (for example, recorded on their original birth certificate).

### Question example:

Sex at birth	Gender
<p><b>What was your sex at birth?</b> (for example, what was recorded on your birth certificate)</p> <p><input type="checkbox"/> male</p> <p><input type="checkbox"/> female</p>	<p><b>What is your gender?</b></p> <p><input type="checkbox"/> male</p> <p><input type="checkbox"/> female</p> <p><input type="checkbox"/> another gender</p> <p>Please state: _____</p>

For older datasets and reports, we need to remember that worldwide, there are limited and inconsistent practices when collecting sex and gender data in ways that reflect the diversity of the population.

This means that many of the existing datasets that we will use, have used the terms "gender" or "sex" interchangeably, and we need to take care when using these, as respondents are more likely to have recorded their gender, or may have recorded their sex at birth due to concerns about discrimination and society views.

<sup>2</sup>

<https://www.stats.govt.nz/consultations/sex-and-gender-identity-statistical-standards-consultation#summary>



# Problem



This section focuses on defining the investigation into the **Population**

© Liz Sneddon

A comparison question needs:

- Categorical variable,
- Numerical variable,
- The words "difference between median",
- Population (**ALL**).

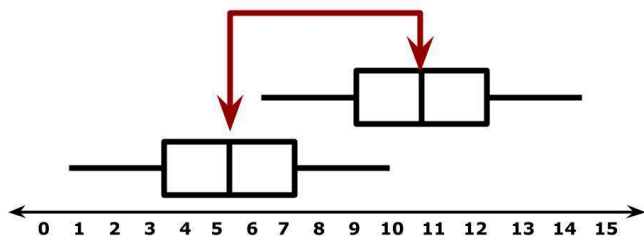
© Liz Sneddon

**FORMAL** investigation question

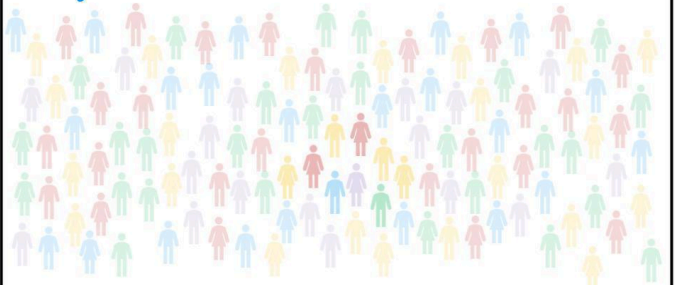
This means that we **DON'T** have a **direction**, but instead write a question asking **WHAT** the **difference between the medians** is.

© Liz Sneddon

Focus on the difference between the medians



**Population:** all the individual members or items that make up a group.



**Notes:**

- You want to compare the centre of one group with the centre of another group. This means you are not expecting that **ALL** the items/people in one group will be bigger or smaller than **ALL** the items/people in the second group. However, you are wanting to compare if **MANY** of the items/people in one group are bigger or smaller than **MANY** of the items/people in the second group.
- The population refers to **EVERYONE** in a group, this is why we suggest you use the word **ALL** when referring to the population. But be careful not to say **ALL people in one group** being bigger or smaller than **ALL people in the second group**, as we are comparing the central tendency of **ONE** population.
- Be careful not to refer to "a" person/item, we are comparing groups of people/items.
- When you refer to the median, it is always the **median of the Numeric variable**, NOT the median of the categorical variable. E.g., **median height**, NOT median boy or median girl.
- The step up from Level 2 Statistics is that instead of just focusing on which median is bigger or smaller, we now accept that there will be some difference between the medians, so at Level 3 the question now shifts to **what the difference** is and we will analyse **the size of the difference**. This means that the investigation question is always "what is the difference between the medians ...".

**Example:**

---

What is the difference between the median weight (kg) of forward and back rugby players from ALL top players listed on the website <http://www.rugby-sidestep-central.com>.

**Exercise 2:**

---

1) In the example given above, identify the following:

<b>Categorical variable</b>	
<b>Two groups being compared</b>	
<b>Numerical variable and units</b>	
<b>Population definition</b>	

2) The questions below all are missing one or more requirements. Identify each of the requirements, and which parts are missing. Then rewrite the question.

a) Are 8-year-old boys generally taller (cm) than 8-year-old girls in NZ?

<b>Categorical variable</b>	
<b>Numerical variable and units</b>	
<b>Difference between median</b>	Yes / No
<b>Population definition</b>	
<b>Rewritten question:</b>	

b) Do all 18-year-old males tend to have a longer right foot than all 18-year-old females in NZ?

<b>Categorical variable</b>	
<b>Numerical variable and units</b>	
<b>Difference between median</b>	Yes / No
<b>Population definition</b>	
<b>Rewritten question:</b>	

- c) I wonder if there is a difference between the average school bag weight for girls and boys?

<b>Categorical variable</b>	
<b>Numerical variable and units</b>	
<b>Difference between median</b>	Yes / No
<b>Population definition</b>	
<b>Rewritten question:</b>	

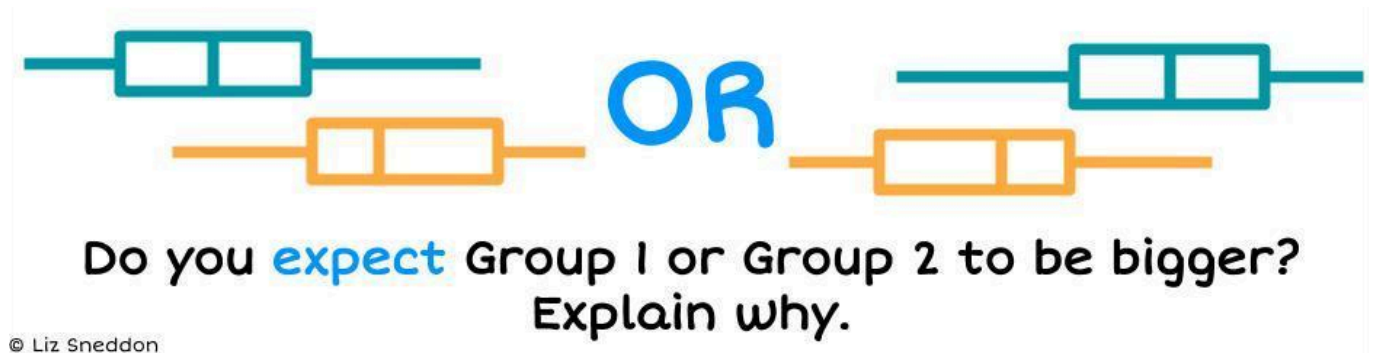
- d) How does the number of text messages all teenage girls send daily compare with the number of text messages all teenage boys send daily, in Auckland?

<b>Categorical variable</b>	
<b>Numerical variable and units</b>	
<b>Difference between median</b>	Yes / No
<b>Population definition</b>	
<b>Rewritten question:</b>	





In the **Problem** section, you want to add a hypothesis about what you might see **BEFORE** you explore the data. You then want to use **research** to support this hypothesis and explain **WHY** there may be a difference between the two groups that you are comparing. Later in both your **Analysis** and **Conclusion** sections you can refer to this research and compare if the data agrees with your hypothesis or not.



## Example:

### Problem:

I wonder if the median weight of teenagers is heavier than the median weight of young children, for all children in NZ.

### Hypothesis:

I expect that the median weight of teenagers would be heavier than young children, as they tend to be taller due to growth spurts.

Research from Kids Health<sup>3</sup> suggests that a major growth spurt usually happens for boys between 10 and 15 years of age, and for girls between 8 and 13 years old. As a child has a growth spurt, they grow taller. As they grow taller the weight of their bones, muscles and organs gets heavier. Therefore, suggesting that teenagers would be heavier than young children.

Now I will investigate the dataset to see if the data supports this hypothesis.

*(Notice that I added a footnote with the website that I got information from).*



<sup>3</sup> <https://kidshealth.org/en/parents/childs-growth.html>



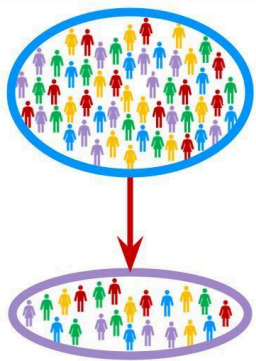

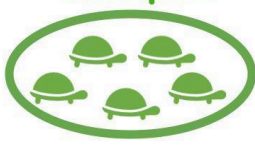

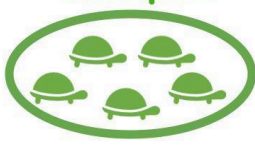

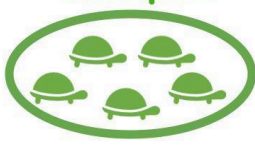




# Plan



One of the most important components of being able to form a conclusion is that the people (or objects) that the data is collected from are **randomly selected**, so that the **sample data is representative of the population**.

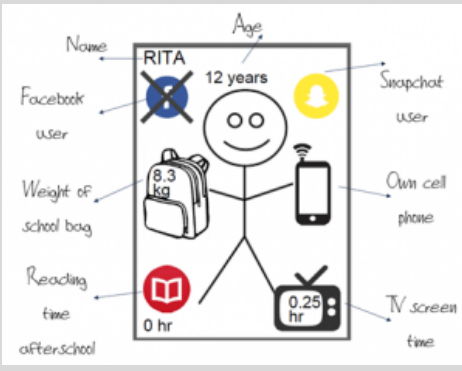
<h3>Sampling method</h3> <p><b>HOW</b> you take a sample from the population</p>  <p>© Liz Sneddon</p>	<table border="0"><tr><td data-bbox="815 1361 1133 1742"><h3>Random sample</h3><p><b>Representative data</b></p><ul style="list-style-type: none"><li>• mixture of characteristics</li><li>• Everyone has the same chance of being selected</li></ul><p>© Liz Sneddon</p></td><td data-bbox="1133 1361 1498 1742"><h3>Biased sample</h3><p><b>NOT representative</b></p><ul style="list-style-type: none"><li>• People have different chances of being selected</li></ul></td></tr></table>	<h3>Random sample</h3>  <p><b>Representative data</b></p> <ul style="list-style-type: none"><li>• mixture of characteristics</li><li>• Everyone has the same chance of being selected</li></ul> <p>© Liz Sneddon</p>	<h3>Biased sample</h3>  <p><b>NOT representative</b></p> <ul style="list-style-type: none"><li>• People have different chances of being selected</li></ul>
<h3>Random sample</h3>  <p><b>Representative data</b></p> <ul style="list-style-type: none"><li>• mixture of characteristics</li><li>• Everyone has the same chance of being selected</li></ul> <p>© Liz Sneddon</p>	<h3>Biased sample</h3>  <p><b>NOT representative</b></p> <ul style="list-style-type: none"><li>• People have different chances of being selected</li></ul>		

## Example

If I do a questionnaire with **only blue-eyed** students, then I have a **biased sample**. This means I do not have any information about people with other coloured eyes (E.g., brown, green, grey, etc.), so my data does not represent the population of all people, only the people with blue eyes.

## Exercise 4:

- 1) Collect a bag of data<sup>4</sup> from your teacher. Take a random sample of 10 students. Record your data below.

 <p><b>Name</b></p>	<b>Age (years)</b>	<b>Do you have Facebook ? Yes / No</b>	<b>Do you have Snapchat ? Yes / No</b>	<b>School bag weight (kg)</b>	<b>Do you have a Cell phone? Yes / No</b>	<b>Reading time yesterday (hours)</b>	<b>TV time yesterday (hours)</b>

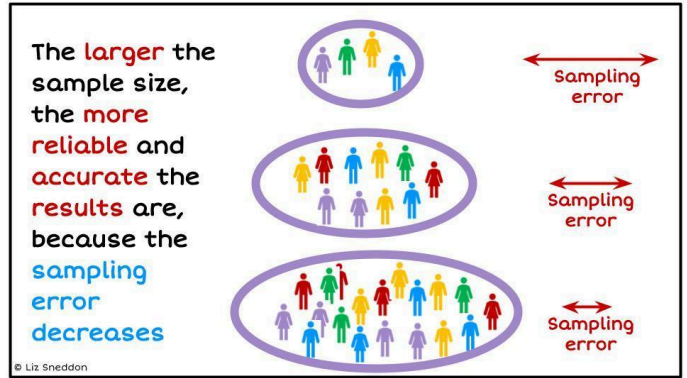
<sup>4</sup> Created by Anna Ferguson: <http://teaching.statistics-is-awesome.org/census-at-school-stick-people-data-cards/>

# Sample size

We want to take a big enough **sample size**, so that the results are **reliable and precise** enough to represent the population.

The more data we have, the greater the precision of our results, and the lower the variation.

With a small sample size, it is much harder to find differences. With a larger sample size, you can find differences more easily.



## Exercise 5:

1) Circle the words that complete the sentences below.

a) Smaller sample sizes take a **shorter / longer** time to collect data, but results are **more / less** precise.

b) Larger sample sizes take a **shorter / longer** time to collect data and results are **more / less** precise.

2) Mrs Sneddon is going to survey 35 girls and 40 boys to investigate their use of iPads at home. Explain why it is ok for the sample sizes to be different, and how these sample sizes affect the variation.

---

---

---

---

---

---

---

---

---

---

# Data / Raraunga



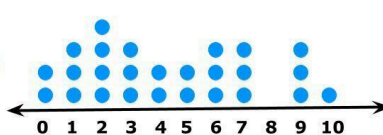
This section focuses processing data from the **Sample**

© Liz Sneddon

Once you have selected the data using a random sampling method, the next step is to process the data, drawing graphs and getting summary statistics.

## Dot plot

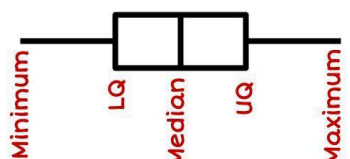
A graph of the data, where each dot is one data value.



© Liz Sneddon

## Box plot

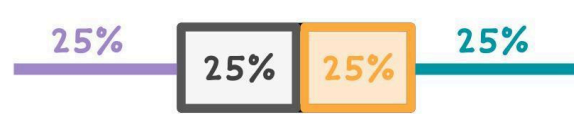
A summary graph of the data.



© Liz Sneddon

## Box plot

Each section of the graph contains 25% of the data.

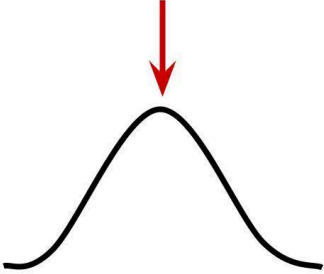


© Liz Sneddon

# Summary Statistics

## Measures of centre

**Centre**



- Mean
- Median
- Mode

© Liz Sneddon

**Mean** =  $\frac{\text{Sum of data}}{\text{Sample size}}$

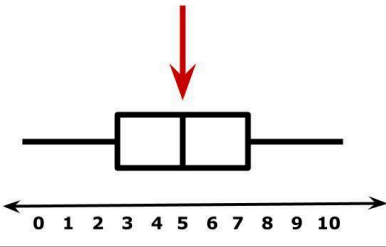
**Median** = middle value

**Mode** = most frequent value

© Liz Sneddon

**Median**

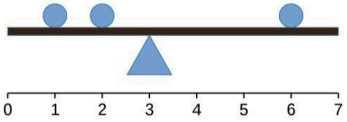
A value where **50%** of the data lies **above & below** it.



© Liz Sneddon

**Mean**

The **average** of a set of values.



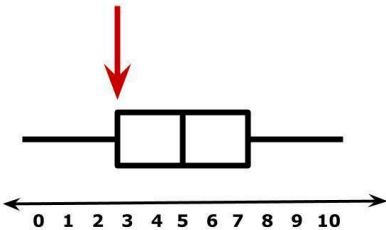
© Liz Sneddon

For numerical data, the mean and median are the preferred measures of centre. For categorical data, the mode is the preferred measure of centre.

## Measures of spread

**Lower Quartile (LQ)**

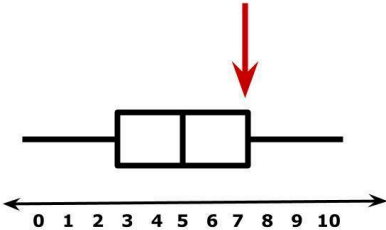
A value where **25%** of the data lies **below** it.



© Liz Sneddon

**Upper Quartile (UQ)**

A value where **25%** of the data lies **above** it.

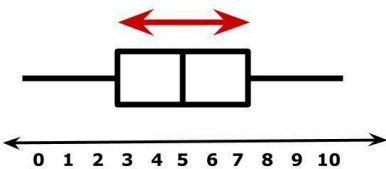


© Liz Sneddon

**Interquartile range**

**IQR = UQ - LQ**

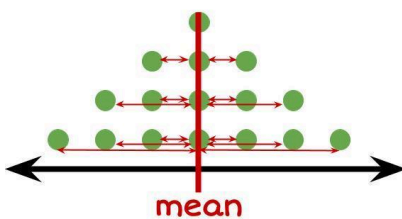
The spread of the **middle 50%** of the data.



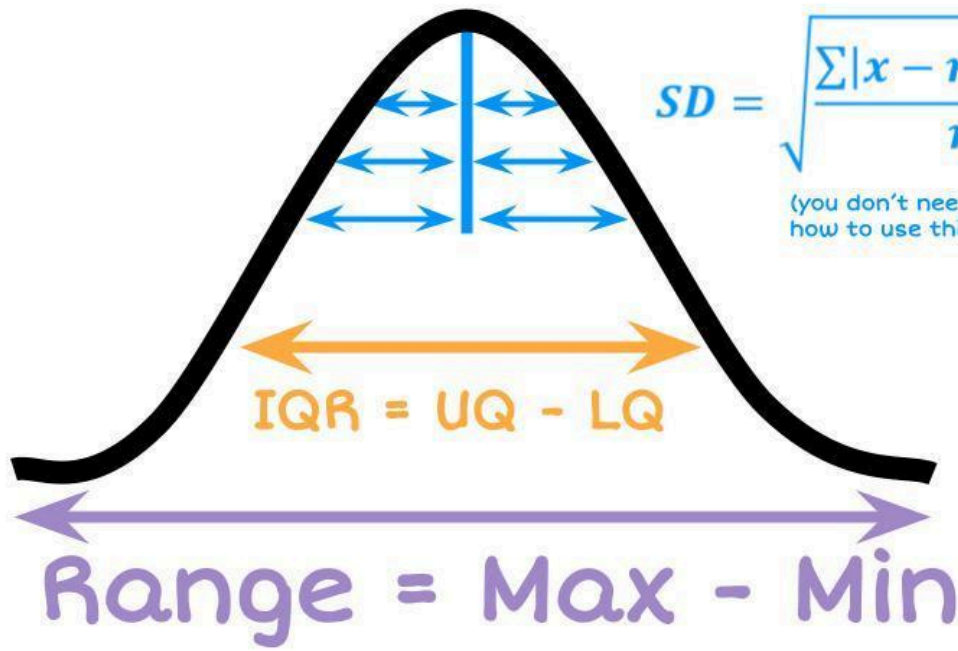
© Liz Sneddon

**Standard deviation**

This measures how far away the data values are from the mean.



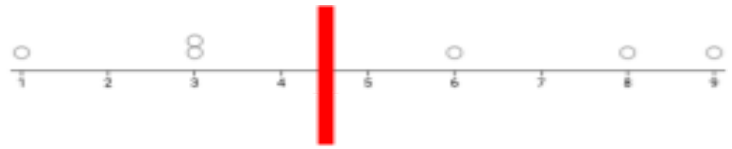
© Liz Sneddon



© Liz Sneddon

## Example:

Estimate the center, and find the mean, median and mode. Then calculate the IQR.

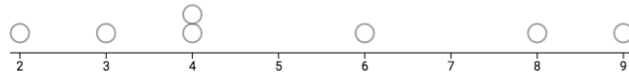


<b>Data</b>	9, 3, 1, 8, 3, 6 Put the numbers in order: 1, 3, 3, 6, 8, 9
<b>Minimum</b>	1
<b>Maximum</b>	9
<b>Median</b>	Find the number(s) in the middle: 1, 3, 3, 6, 8, 9 Find the median = $\frac{3+6}{2} = 4.5$
<b>Mean</b>	$\frac{9+3+1+8+3+6}{6} = 5$
<b>Mode</b>	3
<b>LQ</b>	Take the numbers <b>below</b> the median: 1, 3, 3 Find the middle of these numbers LQ = 3
<b>UQ</b>	Take the numbers <b>above</b> the median: 6, 8, 9 Find the middle of these numbers UQ = 8
<b>IQR</b>	UQ - LQ = 8 - 3 = 5

## Exercise 6:

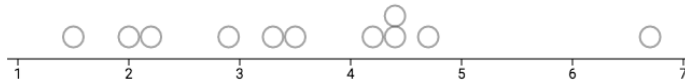
Calculate the summary statistics from the graph, including units.

1)



<b>Data</b>	4, 6, 3, 8, 2, 4, 9
<b>Ordered data</b>	
<b>Minimum</b>	
<b>Maximum</b>	
<b>Median</b>	
<b>Mean</b>	
<b>LQ</b>	
<b>UQ</b>	
<b>IQR</b>	

2)



<b>Data</b>	4.4 4.7 3.5 2.2 4.2 6.7 2.9 4.4 1.5 2.0 3.3
<b>Ordered data</b>	
<b>Minimum</b>	
<b>Maximum</b>	
<b>Median</b>	
<b>Mean</b>	
<b>LQ</b>	
<b>UQ</b>	
<b>IQR</b>	

**Which measures should be used?**

In order of sophistication, our preference for which measure to use (if our data meets certain criteria) is as follows:

Centre:	Spread:
1) Mean 2) Median 3) Mode	1) Standard deviation, 2) IQR, 3) Range.

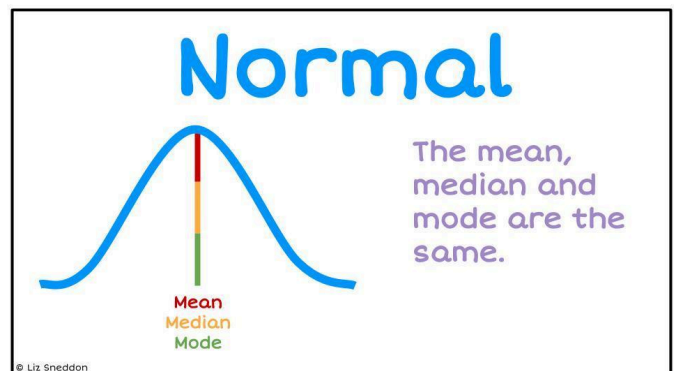
It's important to think about what the best way is to measure the centre and spread, depending on what **shape** the distribution is. Consider how the centre and spread are measured, and which measures of centre and spread will be affected by outliers and skewness, and which measures are the most stable.

We are going to focus on 3 main types of shaped distributions, normal, left skewed, and right skewed.

## For **normally** distributed data:

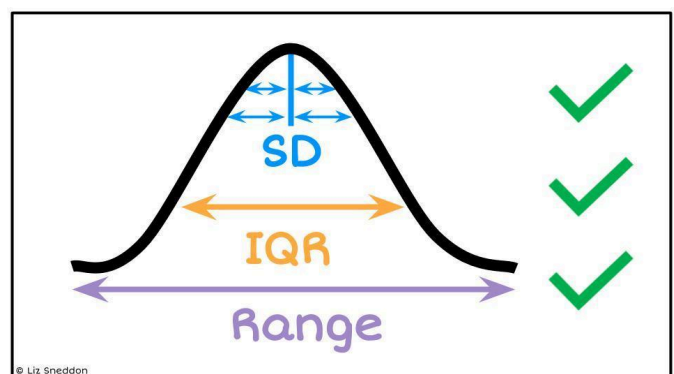
### Centre:

- The mean, median and mode are all **very similar**, therefore they are all **good** measures of centre.
- Therefore, we choose to analyse the **most sophisticated measure**, which is the **mean**.



### Spread:

- The standard deviation, IQR and range are all **good** measures of spread.
- Therefore, we choose to analyse the **most sophisticated measure**, which is the **standard deviation**.

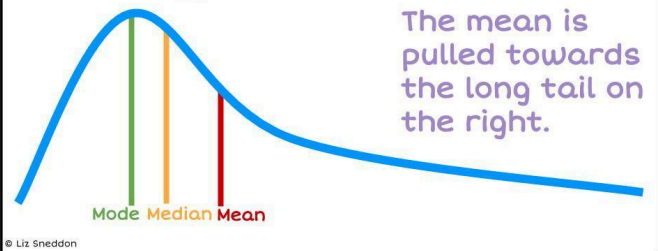


## For **skewed** distributions:

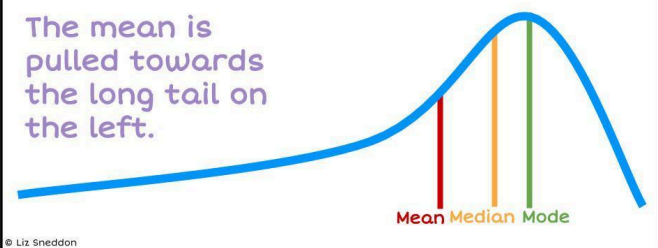
## Centre:

- The mean, median and mode are **different**.
- The mean is pulled in the direction of the longer tail. For right skewed data, the mean is pulled towards the right. For left skewed data, the mean is pulled towards the left.
- The mean is more affected by outliers and skewness, due to the way that it is calculated.
- Therefore, we choose to analyse the 2nd most sophisticated measure, and one that is **NOT** affected by outliers or skewness (due to the way that it is calculated), which is the **median**.

## Right Skew

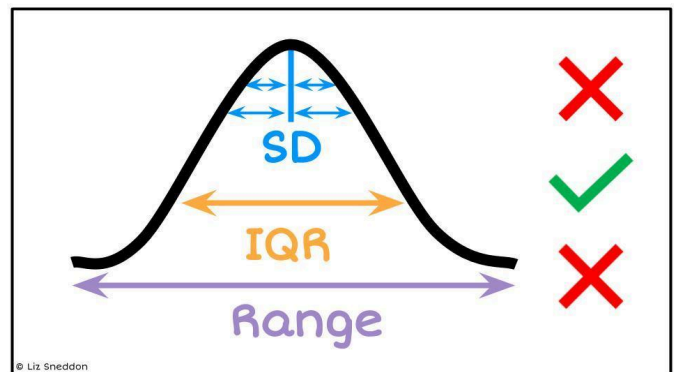


## Left Skew



## Spread:

- The range and the standard deviation are both affected by outliers and skewness, due to the way that it is calculated.
- The IQR is not affected by skewness or outliers, due to the way that it is calculated.
- Therefore, we analyse the one that is **NOT** affected by outliers or skewness, which is the **IQR**.





3) Explain how the range and IQR are calculated.

---

---

---

---

---

---

---

---

4) Explain why the Range is affected by outliers but the IQR is not.

---

---

---

---

---

---

---

---

5) Explain why the IQR is a more stable measure of the spread than the standard deviation if we have outliers.

---

---

---

---

---

---

---

---

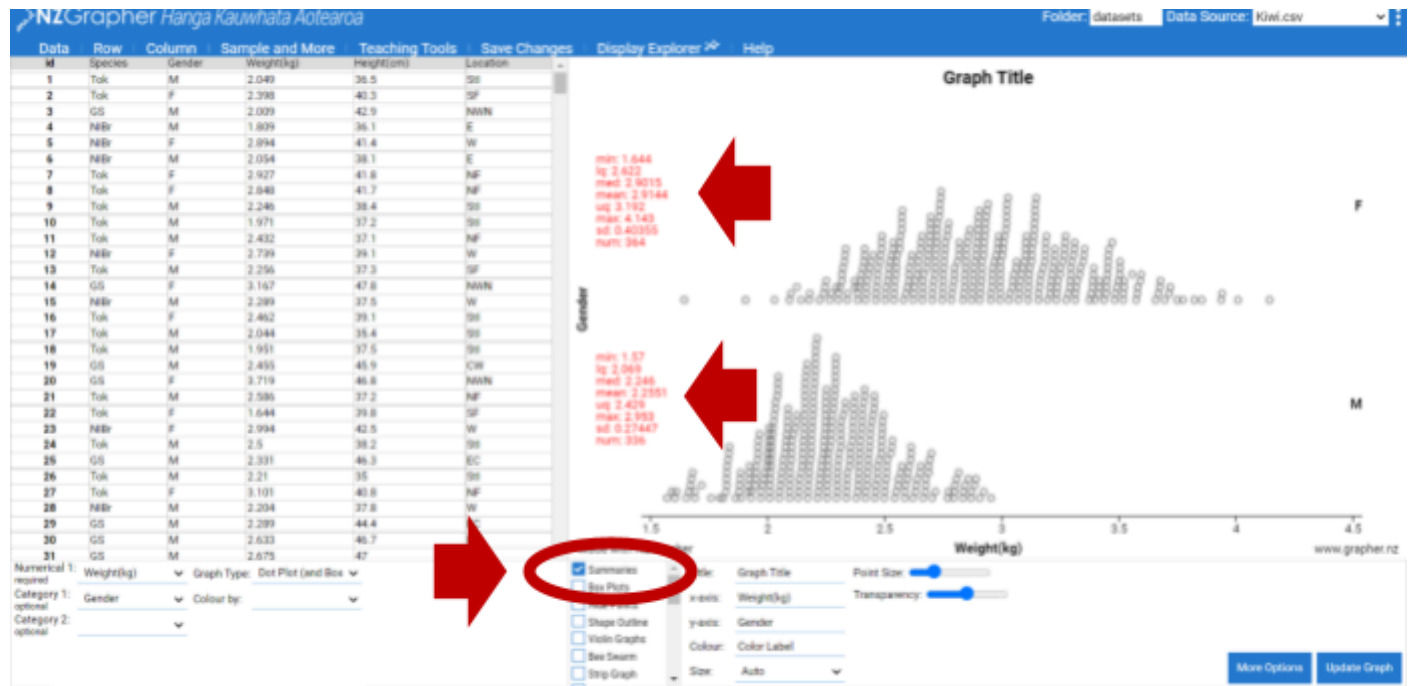




# Dot Plots & Boxplots in NZGrapher

Along with adding the summary statistics, there are several additional options that you need to select when making a Dot plot and Boxplot.

## 1) Add summary statistics to your graph.



Be careful that you don't mix up the **median** and the **mean** values.

min: 1.644  
 lq: 2.622  
 med: 2.9015  
 mean: 2.9144  
 uq: 3.192  
 max: 4.143  
 sd: 0.40355  
 num: 364

med = median  
 mean

## 2) Add titles and axis labels to your graph. Press the Update Graph button when you have finished.



- 3) There are several options for box plots, one that sits **on top of** the data, and one that sits **above** the data (tick the High Box plot option).

The image shows the NZGrapher interface with two types of box plots. The top plot is a 'High Box Plot' where the box is positioned above the data points. The bottom plot is a 'Standard Box Plot' where the box is overlaid on the data points. A red arrow points from the 'High' label to the top plot, and another red arrow points from the 'Standard' label to the bottom plot. Below the plots, the 'Box Plots' settings are visible, with 'Box Plots' and 'High Box Plot' options checked and circled in red.

## Bootstrap Graph in NZGrapher

At Level 3 there is an additional graph that we need to add, and that is the bootstrap graph – median. We will explore what this graph is and how to interpret the results in the Conclusion section.

Change the graph type to “**Bootstrap Confidence Interval - Median**”

The image shows the NZGrapher interface with a data table on the left and a graph on the right. The data table has columns for ID, Species, Gender, Sample, Weighting, Height(cm), and Location. The graph is titled 'Kiwi Investigation' and shows a distribution of weight (kg) for two genders: Female (F) and Male (M). The x-axis is 'Weight (kg)' and the y-axis is 'Gender'. A red arrow points from the 'Bootstrap Confidence Interval - Median' option in the settings to the graph. The graph shows a distribution of weight (kg) for two genders: Female (F) and Male (M). The x-axis is 'Weight (kg)' and the y-axis is 'Gender'. A red arrow points from the 'Bootstrap Confidence Interval - Median' option in the settings to the graph.

# NZGrapher Tips

Here are some additional tips which I find useful.

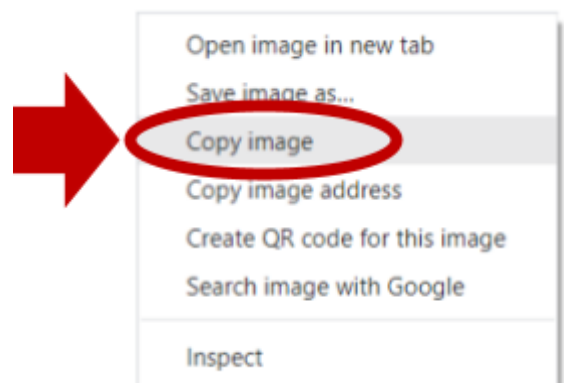
**Tip 1:** You can tidy up the way the dots look, by selecting **Stack Dots**, and you can also change the **point size** to better see the shape of the data.



**Tip 2:** There are some new graphs that you can change / add. I like the Shape outline as the overall shape can be seen more easily.



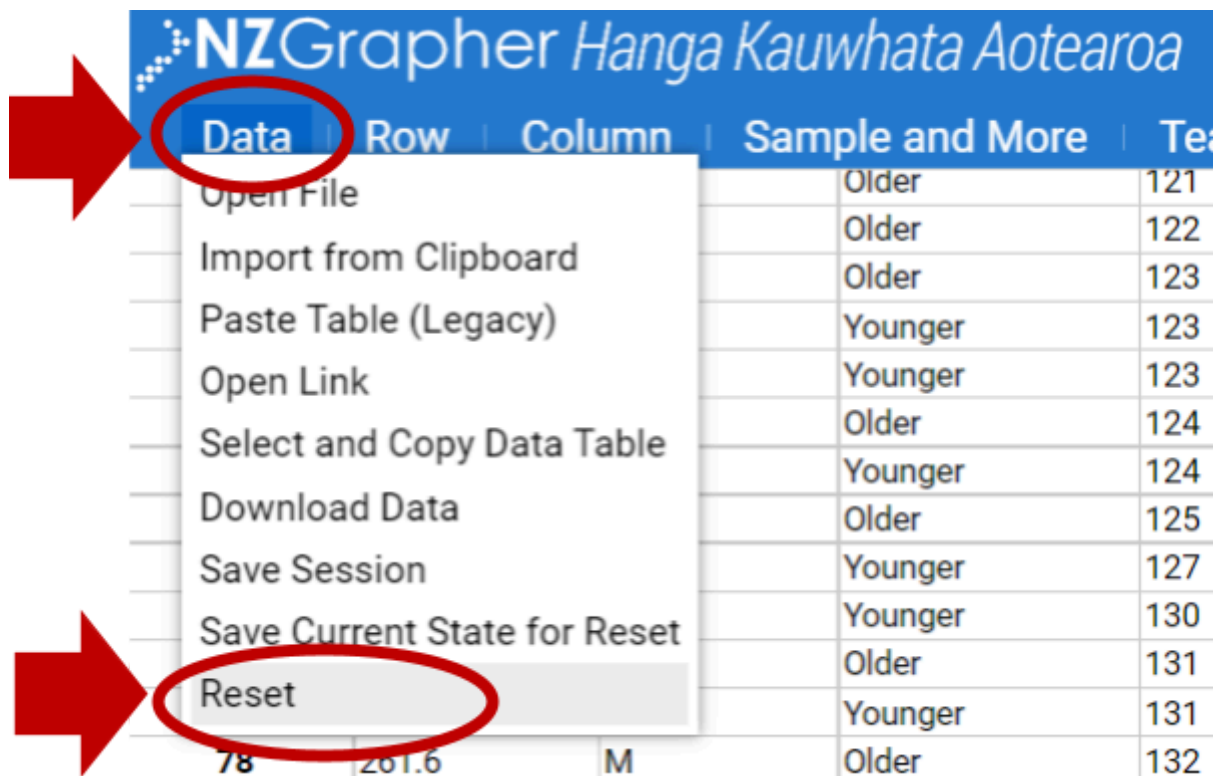
**Tip 3:** You can copy each graph by right clicking over the image, and then selecting **Copy Image**.



**Tip 4:** You can select the **Box plot (no Outlier)** which will draw the box plot **EXCLUDING** the outliers in the tail. This allows you to easily identify the outliers.



**Tip 5:** If **you** make a mistake when you take a random sample, then go to the Data menu and select the **Reset** button. This will reload the full dataset and you can start again.



# Putting NZGrapher together

These are the three graphs I usually produce to help me analyse and draw conclusions using the data.

## Graph 1: Dot Plot – Analysis of features

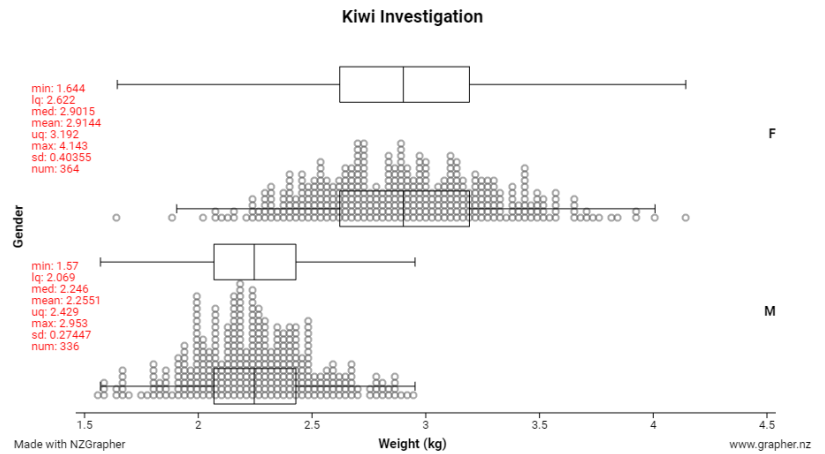
Select the following options:	Graph
<input checked="" type="checkbox"/> Summaries <input type="checkbox"/> Box Plots <input checked="" type="checkbox"/> High Box Plot <input type="checkbox"/> Hide Points <input checked="" type="checkbox"/> Shape Outline <input type="checkbox"/> Violin Graphs <input type="checkbox"/> Bee Swarm <input type="checkbox"/> Strip Graph <input type="checkbox"/> Box (No Whisker) <input type="checkbox"/> Box (No Outlier) <input type="checkbox"/> DBM & OVS <input type="checkbox"/> Numbers ½ ¾ Rule: None <hr/> <input type="checkbox"/> Informal C-I <input type="checkbox"/> C-I Limits <input type="checkbox"/> C-I Highlight <input type="checkbox"/> Point Labels <input type="checkbox"/> Mean Dot <input checked="" type="checkbox"/> Stack Dots <input type="checkbox"/> Gridlines	<p style="text-align: center;"><b>Kiwi Investigation</b></p> <p>min: 1.644            lq: 2.622            med: 2.9015            mean: 2.9144            uq: 3.192            max: 4.143            sd: 0.40355            num: 364</p> <p>min: 1.57            lq: 2.069            med: 2.246            mean: 2.2551            uq: 2.429            max: 2.958            sd: 0.27447            num: 336</p> <p>Gender: F, M</p> <p>Weight (kg)</p> <p>Made with NZGrapher <span style="float: right;">www.grapher.nz</span></p> <p>Don't forget that you can change the point size if you have small sample sizes:</p> <p>Point Size: <input type="range"/></p> <p>Transparency: <input type="range"/></p>

## Graph 2: Dot Plot – Outliers

Select the following options:

## Graph

- Summaries
  - Box Plots
  - High Box Plot
  - Hide Points
  - Shape Outline
  - Violin Graphs
  - Bee Swarm
  - Strip Graph
  - Box (No Whisker)
  - Box (No Outlier)
  - DBM & OVS
  - Numbers
  - $\frac{1}{2}$   $\frac{3}{4}$  Rule:
  - None
- 
- Informal C-I
    - C-I Limits
    - C-I Highlight
  - Point Labels
  - Mean Dot
  - Stack Dots
  - Gridlines



Graph 3: Bootstrap Confidence Interval – Medians

Select the following options:	Graph
<input type="checkbox"/> Strip Graph <input type="checkbox"/> Point Labels <input checked="" type="checkbox"/> Stack Dots <input type="checkbox"/> Gridlines <input type="checkbox"/> Thick Lines <input checked="" type="checkbox"/> Show ID of Removed Points	

### Exercise 8:

Go to NZGrapher and form appropriate graphs for each of the datasets below (don't forget titles, axis labels etc):

- 1) Select the Marathon dataset and use the variables "**Minutes**" and "**AgeGroup**".
- 2) Select the Kiwi dataset and use the variables "**Weight**" and "**Gender**".
- 3) Select the BallWear dataset and use the variables "**Gender**" and "**AmountSpent**".

# Analysis / Tātari

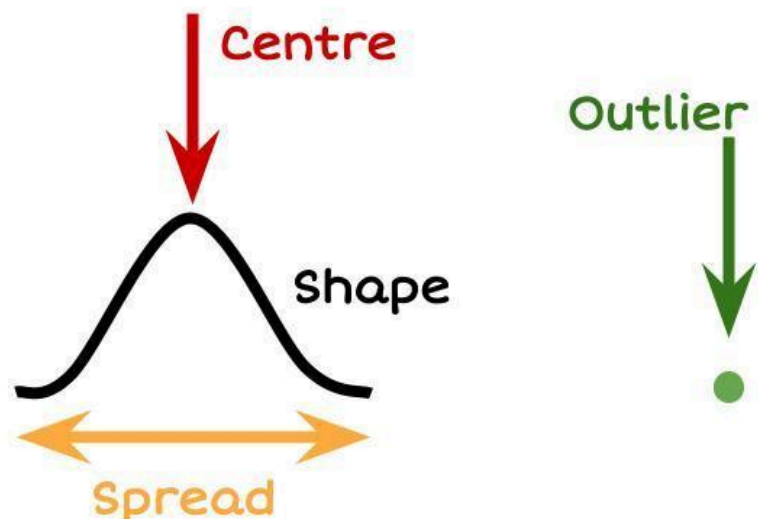


This section focuses analysing the features of the **Sample**

© Liz Sneddon

## Features

- Shape
- Centre
- Spread
- Outliers



© Liz Sneddon

Remember to use **suggestive language** when you analyse the features and patterns in your **sample data**.

We will go through each of these features separately, before putting it altogether.

# Shape

When we analyse the shape, we are taking a step back and thinking about what shape the distribution is that this sample has come from. It's like standing at the back of the classroom and seeing the general pattern of the data on the board.

## Shape

The pattern the data forms when plotted. It varies depending on number of peaks, symmetry and skewness.

Right skewed  
Left skewed  
Normal  
Bimodal  
Uniform

© Liz Sneddon

## Skew

When the data on one side is tightly packed and the other side is much more spread out.

Right skewed  
Left skewed

© Liz Sneddon

Identifying shapes:

<b>Normal distribution</b>	<b>Left skewed</b>	<b>Right Skewed</b>
<b>Bimodal</b>	<b>Uniform</b>	

## Example:

Draw shapes like this:



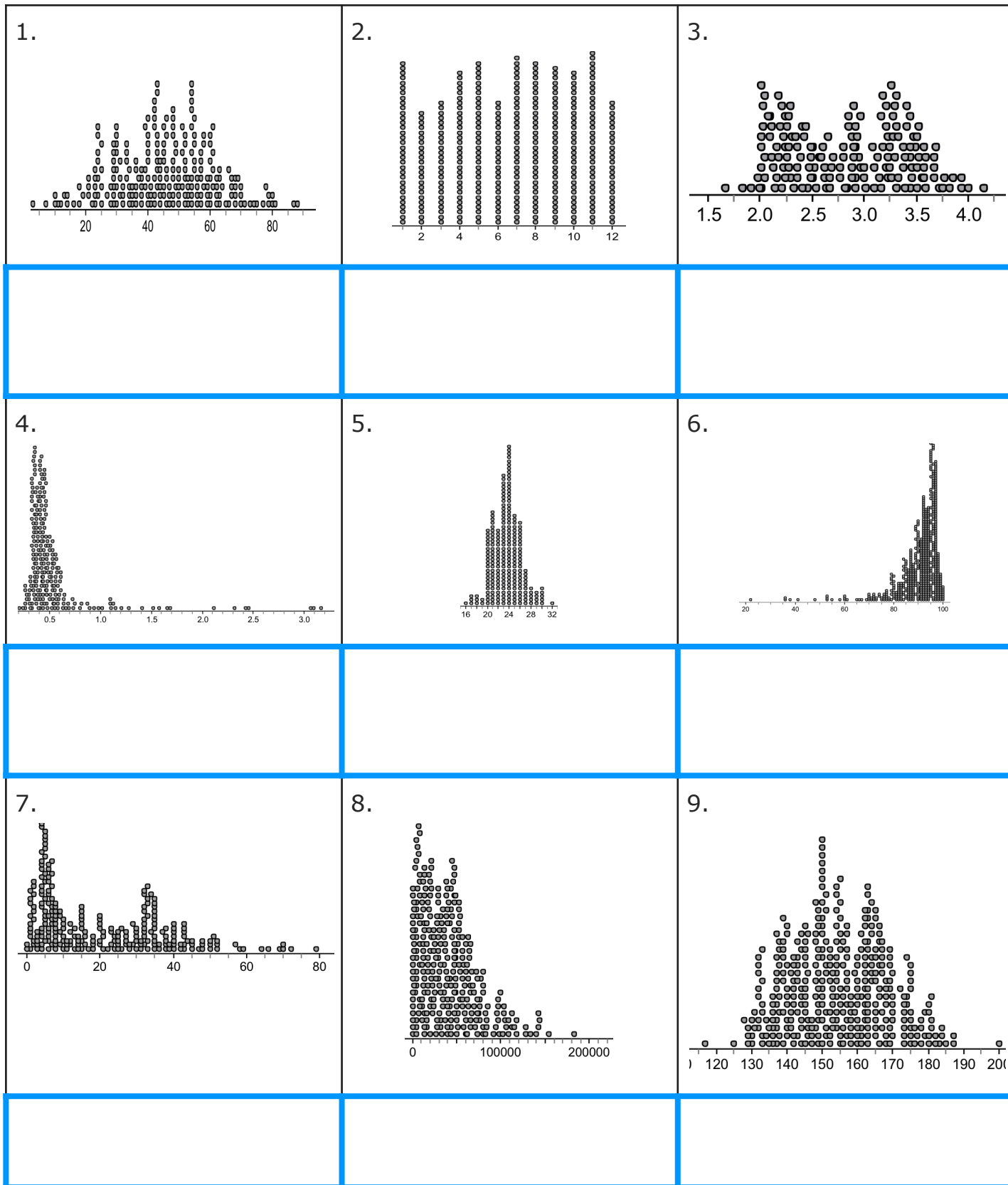
Not like this:



The data has 2 peaks, so looks approximately bimodal in shape.

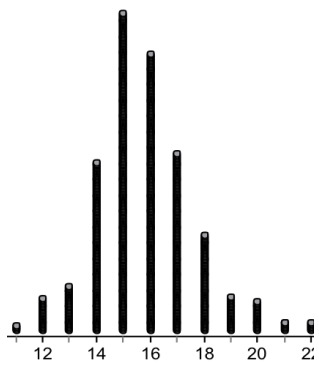
## Exercise 9:<sup>5</sup>

Sketch over the top of each graph and then state what shape it most closely matches.

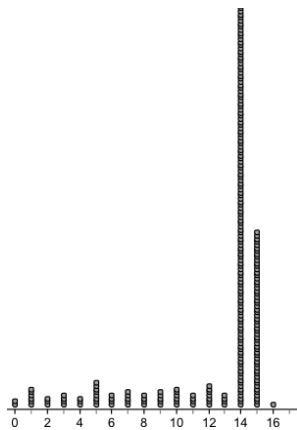


<sup>5</sup> Thanks to Dr Pip Arnold for the graphs.

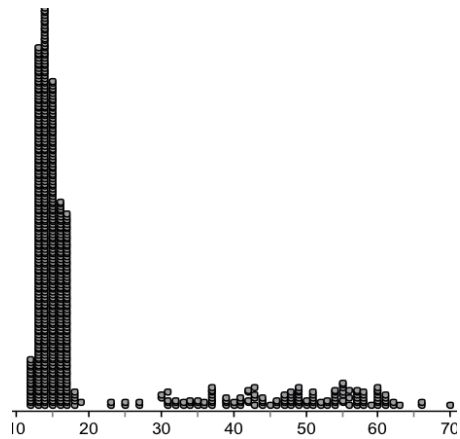
10.



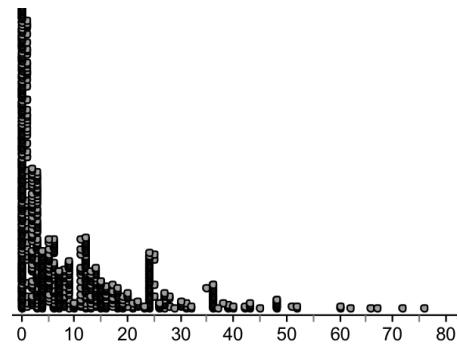
11.



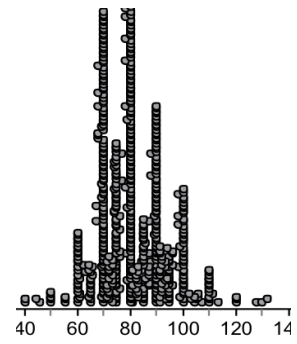
12.



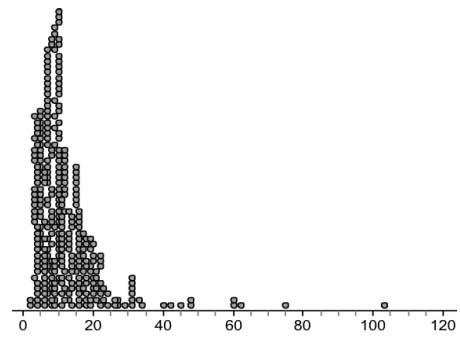
13.



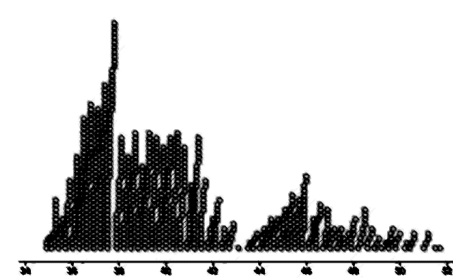
14.



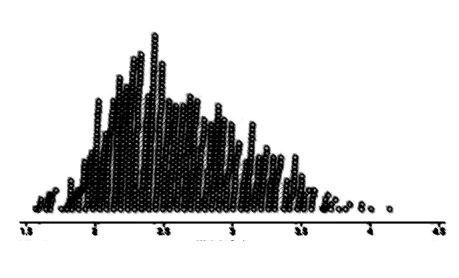
15.



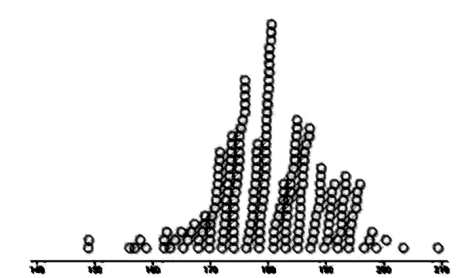
16.



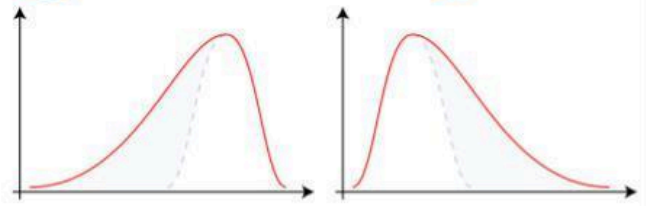
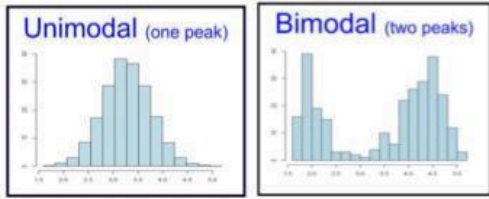
17.



18.



# Justifying Shape



Number of peaks

Skewness

Symmetry



Symmetrical

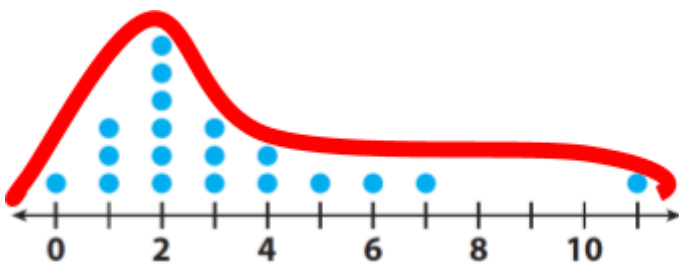
Asymmetrical

© Liz Sneddon

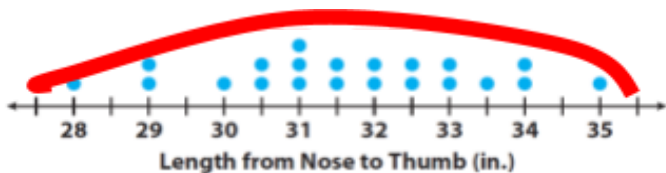
## Examples:



There are two hills / peaks, the data is not symmetric in shape and the tails on each side are different. Overall, the shape of the data appears to be bimodal.



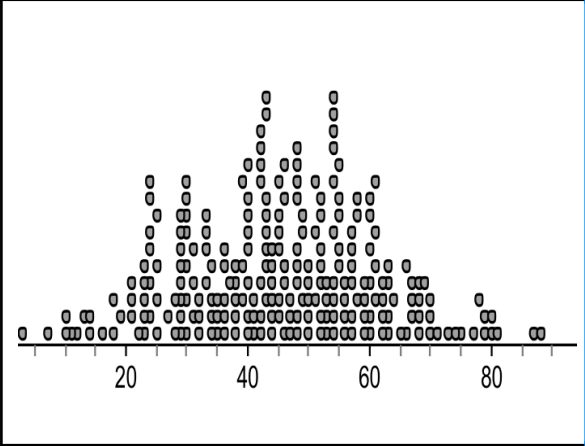
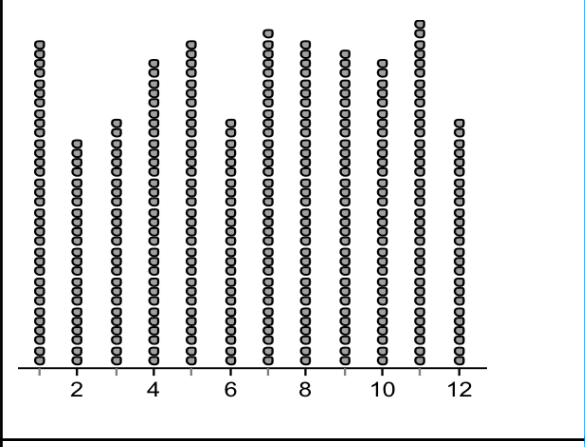
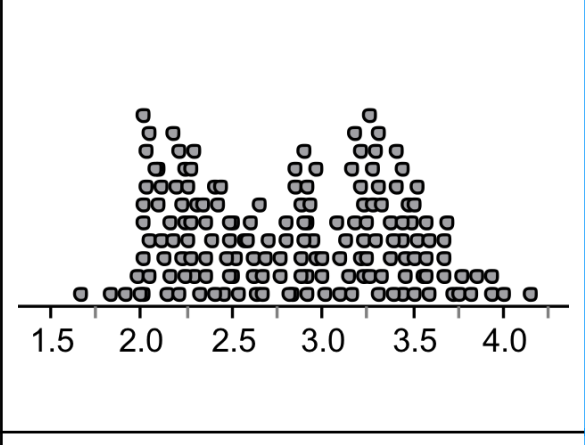
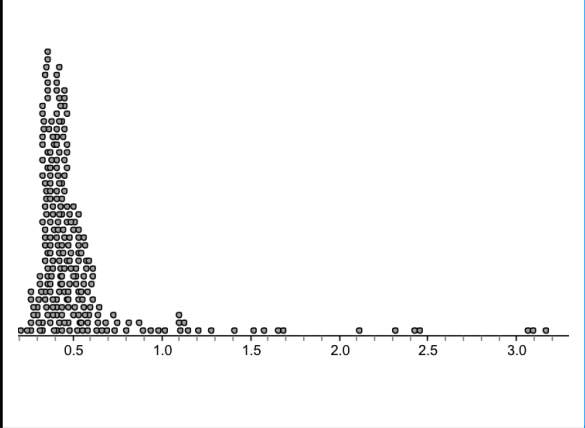
The shape of the data looks right skewed, as there is a longer tail on the right, with one hill / peak. It is not symmetric.

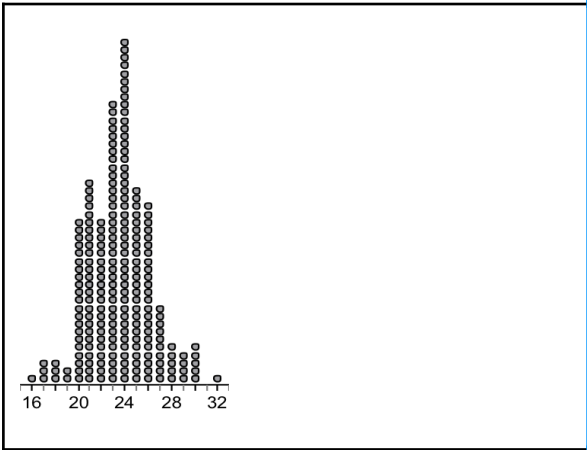


The shape of the data looks approximately normal, as there is one hill / peak, it is roughly symmetric and the tails on both sides are similar.

## Exercise 10:

For each graph, select the correct symmetry, peaks, and tail description.

	<p>Symmetric (circle one):    Yes / No</p> <p>Number of hills / peaks:    0   1   2</p> <p>Bell shaped curve:            Yes / No</p> <p>Tails (tick only if appropriate):</p> <p>    Left hand tail longer</p> <p>    Right hand tail longer</p>
	<p>Symmetric (circle one):    Yes / No</p> <p>Number of hills / peaks:    0   1   2</p> <p>Bell shaped curve:            Yes / No</p> <p>Tails (tick only if appropriate):</p> <p>    Left hand tail longer</p> <p>    Right hand tail longer</p>
	<p>Symmetric (circle one):    Yes / No</p> <p>Number of hills / peaks:    0   1   2</p> <p>Bell shaped curve:            Yes / No</p> <p>Tails (tick only if appropriate):</p> <p>    Left hand tail longer</p> <p>    Right hand tail longer</p>
	<p>Symmetric (circle one):    Yes / No</p> <p>Number of hills / peaks:    0   1   2</p> <p>Bell shaped curve:            Yes / No</p> <p>Tails (tick only if appropriate):</p> <p>    Left hand tail longer</p> <p>    Right hand tail longer</p>



Symmetric (circle one):    Yes / No

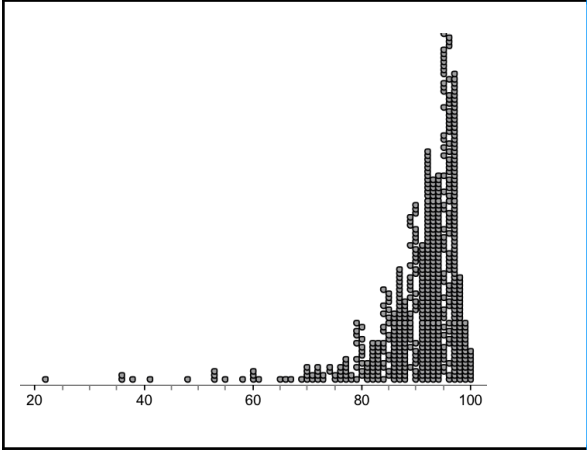
Number of hills / peaks:    0   1   2

Bell shaped curve:            Yes / No

Tails (tick only if appropriate):

    Left hand tail longer

    Right hand tail longer



Symmetric (circle one):    Yes / No

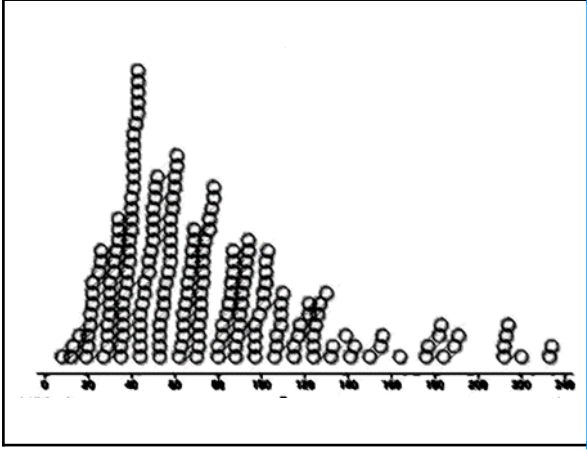
Number of hills / peaks:    0   1   2

Bell shaped curve:            Yes / No

Tails (tick only if appropriate):

    Left hand tail longer

    Right hand tail longer



Symmetric (circle one):    Yes / No

Number of hills / peaks:    0   1   2

Bell shaped curve:            Yes / No

Tails (tick only if appropriate):

    Left hand tail longer

    Right hand tail longer

# Shapes with NZGrapher

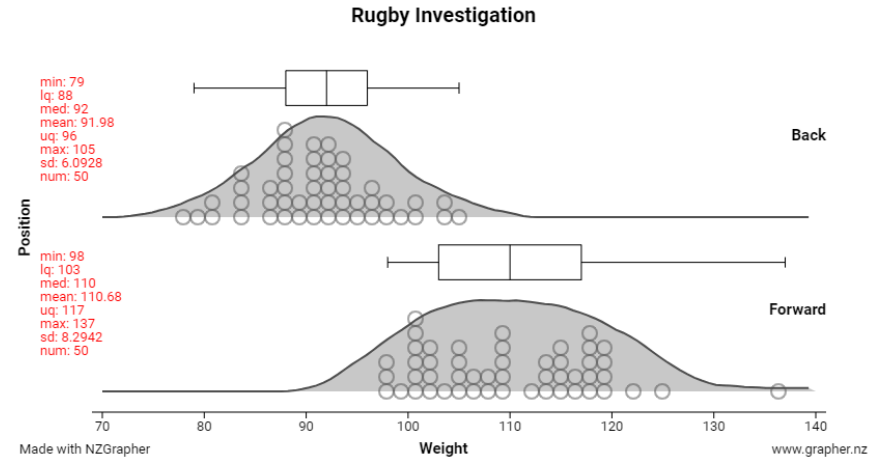
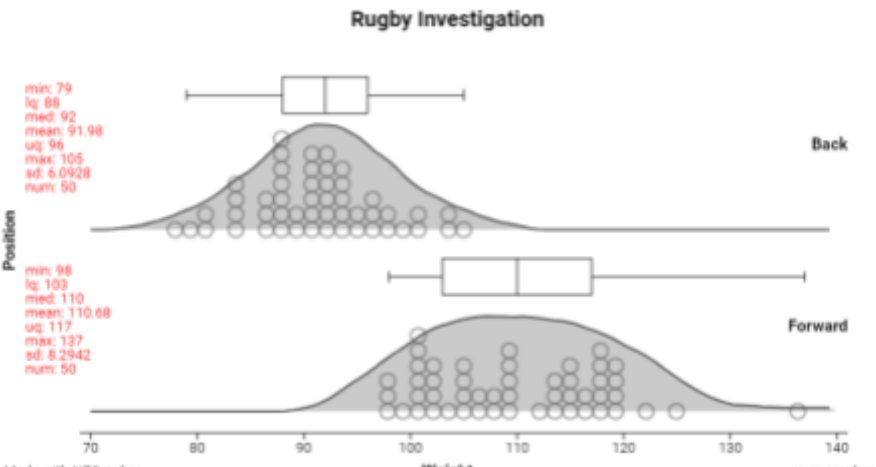
When NZGrapher draws the shape outline over each dataset, it is endeavoring to show what the **shape of the population BEHIND the sample of data** might look like. This is why the curves are smooth, and this is the shape we are trying to identify.

When NZGrapher we need to remember that when we identify the shape, we **ONLY** consider the shape between the **minimum and maximum** values of **EACH group**. (Use the whiskers on the box plot to remind you of these minimums and maximums).

The best way to do this, is to use a highlighter to colour in the data section of the graph to allow you to focus on only that section.

You also need to **IGNORE** any of the tails that extend beyond the data.

## Example:

<p><b>Step 1:</b> Get the graph from NZGrapher</p>	 <p>Rugby Investigation</p> <p>min: 79 Iq: 88 med: 92 mean: 91.98 uq: 96 max: 105 sd: 6.0928 num: 50</p> <p>Back</p> <p>min: 98 Iq: 103 med: 110 mean: 110.68 uq: 117 max: 137 sd: 8.2942 num: 50</p> <p>Forward</p> <p>Position</p> <p>Weight</p> <p>Made with NZGrapher</p> <p>www.grapher.nz</p>
<p><b>Step 2:</b> Highlight the areas between the minimum and maximum of each group, ignoring the tails.</p>	 <p>Rugby Investigation</p> <p>min: 79 Iq: 88 med: 92 mean: 91.98 uq: 96 max: 105 sd: 6.0928 num: 50</p> <p>Back</p> <p>min: 98 Iq: 103 med: 110 mean: 110.68 uq: 117 max: 137 sd: 8.2942 num: 50</p> <p>Forward</p> <p>Position</p> <p>Weight</p> <p>Made with NZGrapher</p> <p>www.grapher.nz</p>
<p><b>Step 3:</b> Identify the shape in context, referring to both the numeric and categorical variables. Then justify both shapes.</p>	<p>The shape of weights for the sample of back rugby players is approximately Normal, because the weights are roughly symmetrical, unimodal and follow a bell-shaped curve.</p> <p>The shape of weights for the forwards rugby players is skewed to the right, because the weights are unimodal, asymmetric and more spread out the right-hand side.</p>



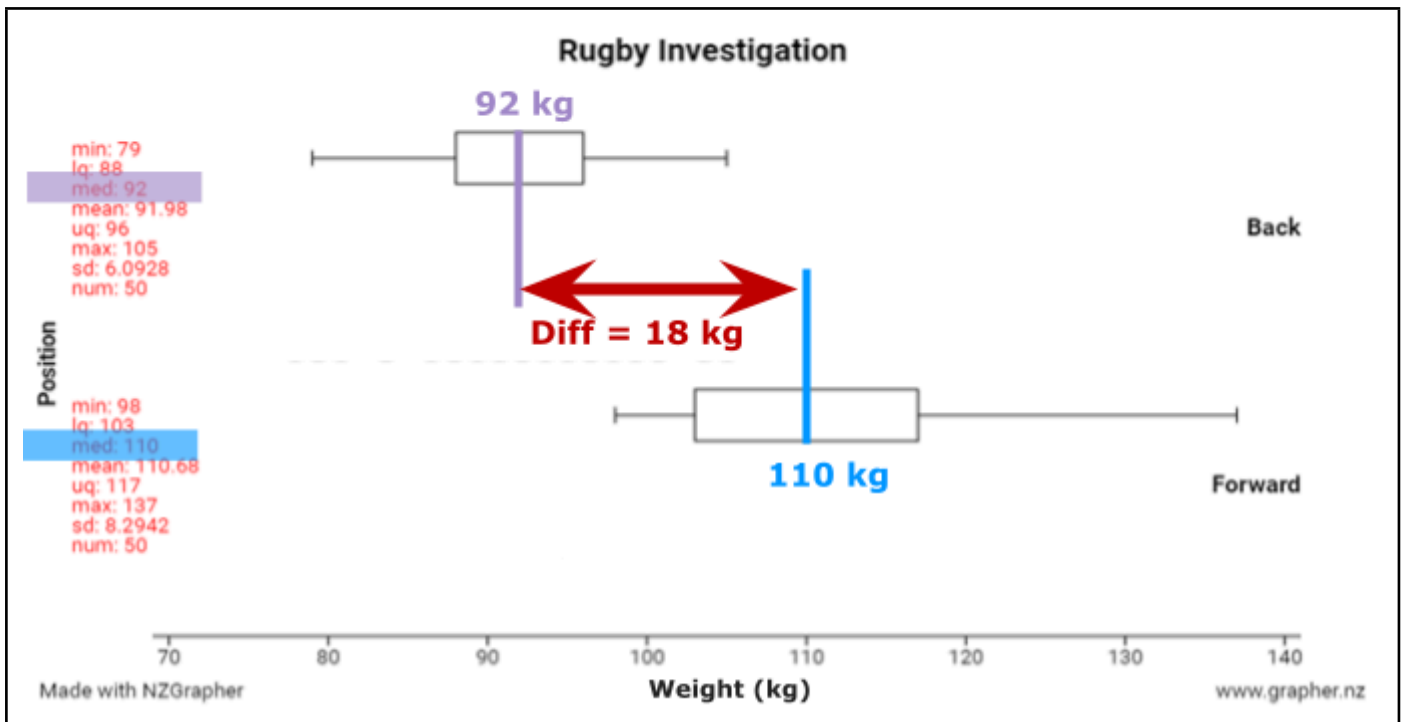






For Excellence, tell the story and connect to research, explaining why one group may (or may not) be bigger or smaller than the second group.

### Example:



The median weight for my sample of backs is 92kg. The median weight for my sample of forwards is 110 kg.

$$\begin{aligned} \text{Difference} &= \text{forwards median} - \text{backs median} \\ &= 110 \text{ kg} - 92 \text{ kg} \\ &= 18 \text{ kg} \end{aligned}$$

In my sample, the median weight for forwards rugby players is heavier than backs by 18 kg, as forwards need more muscles and weight to be able to both hold the line and push the line forwards, whereas backs tend to need to run fast, which often is less bulk than rugby forwards.

### Exercise 12:

Compare the medians for the graphs below.

1)









# Comparing Spread

## Spread

The spread in data is the measure of how far the numbers in a data set are away from the centre.

© Liz Sneddon

## Spread: Compare the IQR's

© Liz Sneddon

## Comparing Spread

When comparing spread, we look for **large** differences, e.g. **double** the width, **triple** the width etc.

© Liz Sneddon

## Interquartile range

### IQR = UQ - LQ

The spread of the **middle 50%** of the data.

© Liz Sneddon

The IQR is the middle 50% of your sample data, and in our boxplot, it is the box.

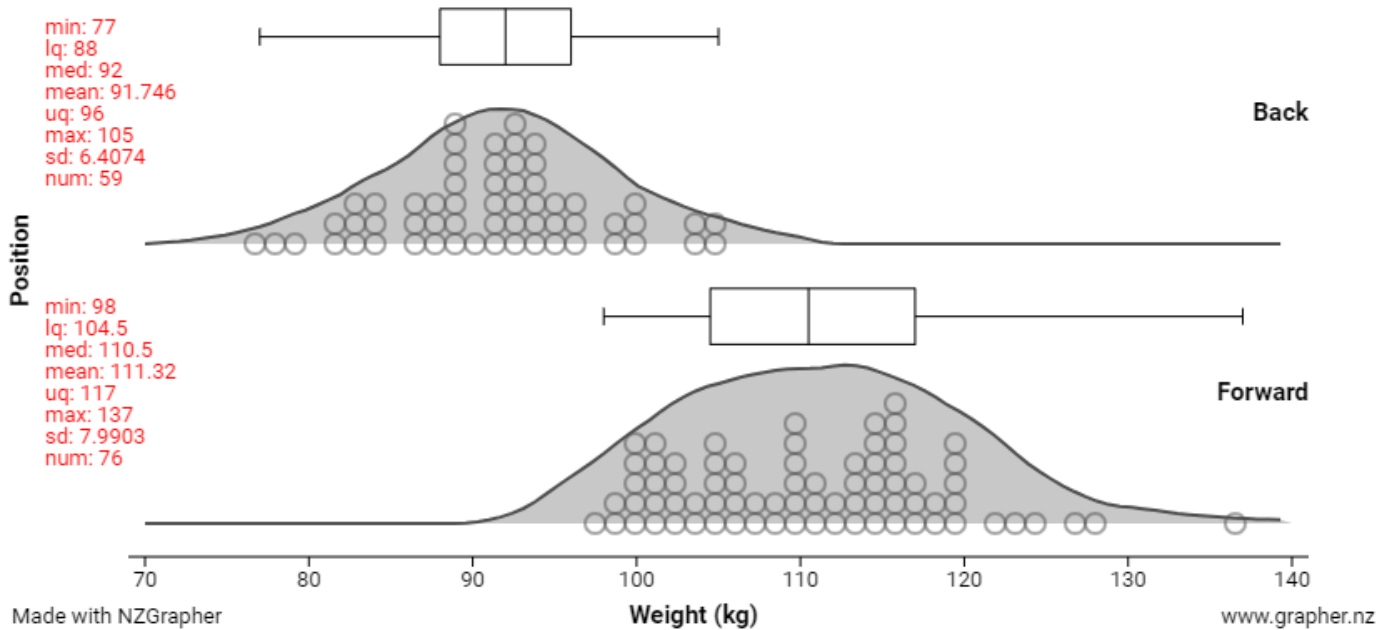
When comparing the spread, don't calculate the difference between the IQR's, look for whether **one IQR is a multiple of the other**. E.g., double, triple, etc.

For Merit, make sure that you include the IQR values and units of each group as evidence.

For Excellence, you **MAY** be able to tell the story, explaining why one group may (or may not) be more or less spread out than the second group. It is not always easy to explain why one group is more or less spread out, so be careful not to confuse this with comparing the centres and medians.

## Example:

### Rugby Investigation



$$\begin{aligned} \text{IQR (backs)} &= \text{UQ} - \text{LQ} \\ &= 96 - 88 \\ &= 8 \text{ kg} \end{aligned}$$

The IQR of weights for back rugby players is 8kg.

$$\begin{aligned} \text{IQR (forwards)} &= \text{UQ} - \text{LQ} \\ &= 117 - 110.5 \\ &= 6.5 \text{ kg} \end{aligned}$$

The IQR of weights for forwards rugby players is 6.5kg.

In the sample, the spread of the middle 50% of weights of back rugby players is a **little wider** than the spread of the middle 50% of weights of forward rugby players.









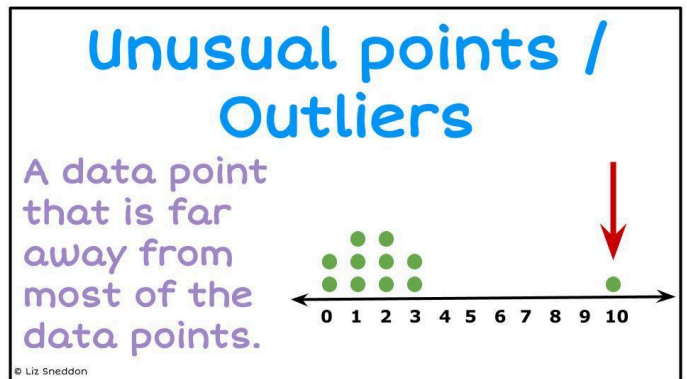
---

---

# Outliers

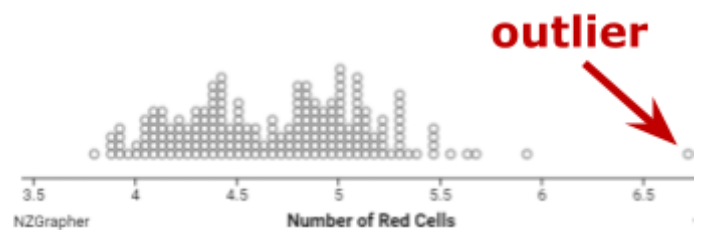
---

When we look for outliers, you are looking for data that is a long way away from the rest of the data. For example, if a student with a height of 2 meters entered the classroom, their height would be unusual. Or if a person who was 100 years old entered the classroom, that would be unusual. This is not one of the main patterns we look for, but if you notice an outlier, identify it.

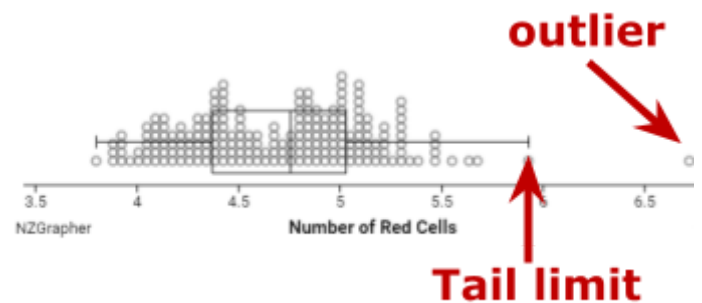


## Example 1:

First look at the graph and identify any points that are a **long way away**.



Then confirm it by adding the **Box plot (no outlier)**. Notice that the whisker stops at the value 5.8, so any data points after this are outliers.



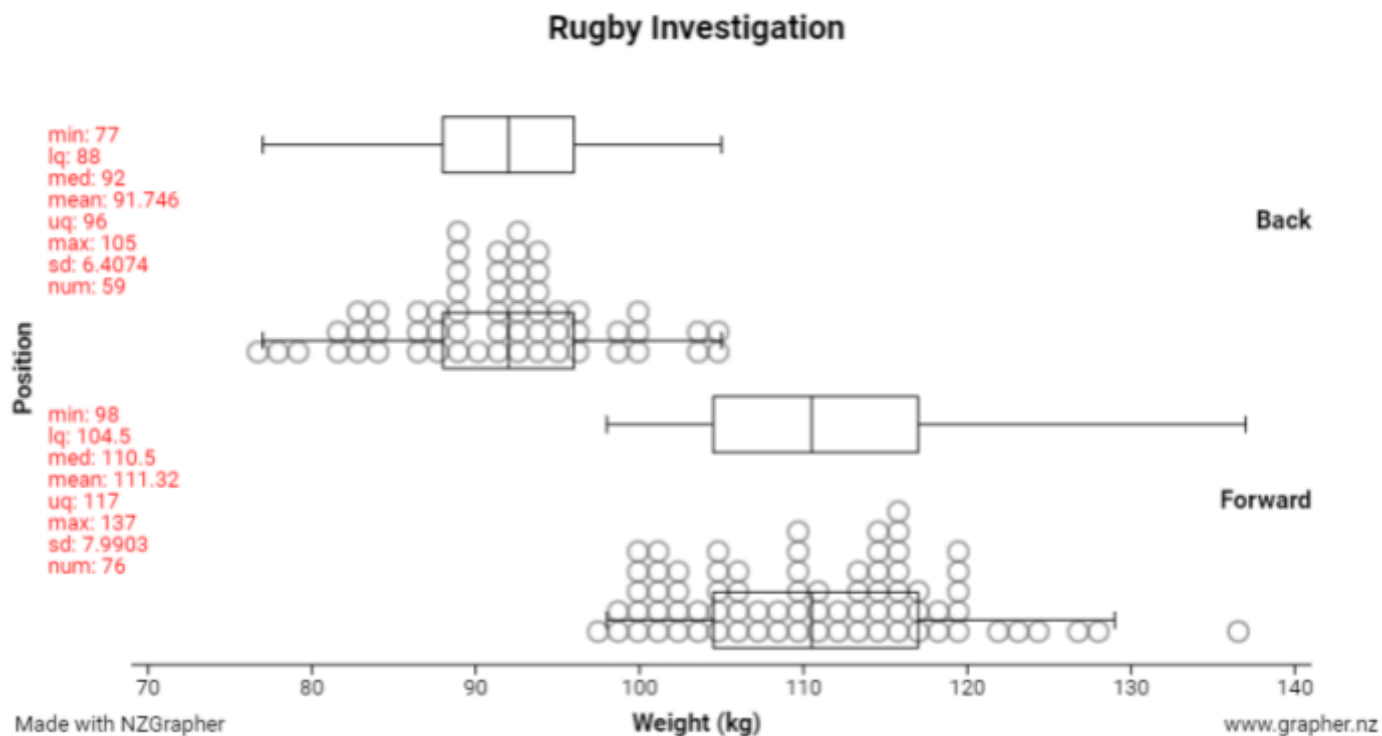
There is an outlier, a person with around 6.8 red blood cells per microliter.

## Example 2:

In the graph below, you can see 2 box plots for each group. The high box plot covers the weights from the minimum to the maximum value. The box plot that is overlaid over the data points, is the Box plot (No outliers).

The key is to compare the ends of the Box plot (No outliers) and see if there are any data points (weights) that are outside of the whiskers.

I have highlighted the 2 outliers in the forwards. Notice that they are outside the ends of the Box plot (No Outliers) – but the high box plot still extends to the minimum and maximum values.



There are two outliers, who both play in a forwards position in rugby. One has a smaller than expected weight of around 96 kg, and another with a weight higher than expected at around 136 kg.









---

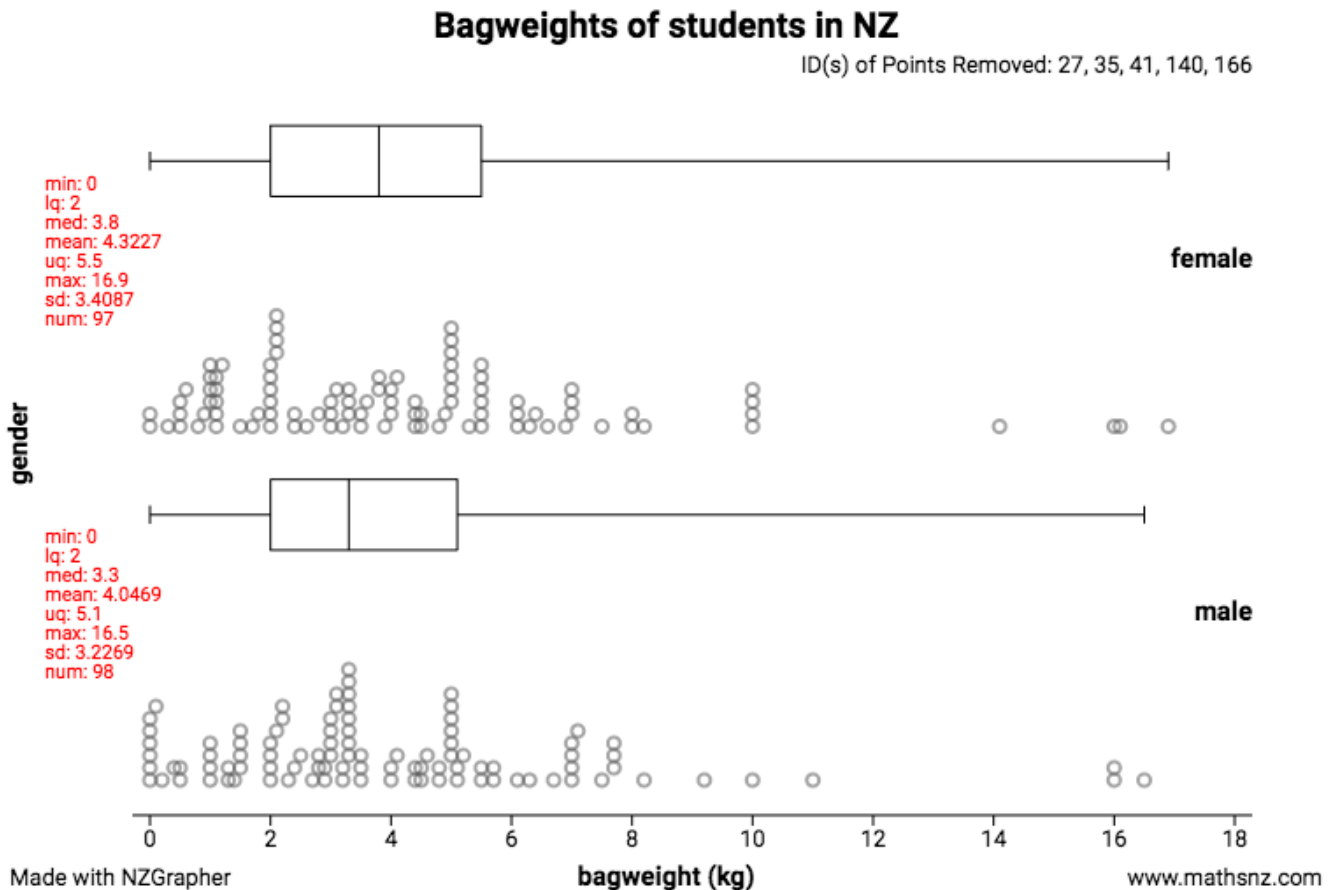
---

# Full Analysis

---

## Example:

---

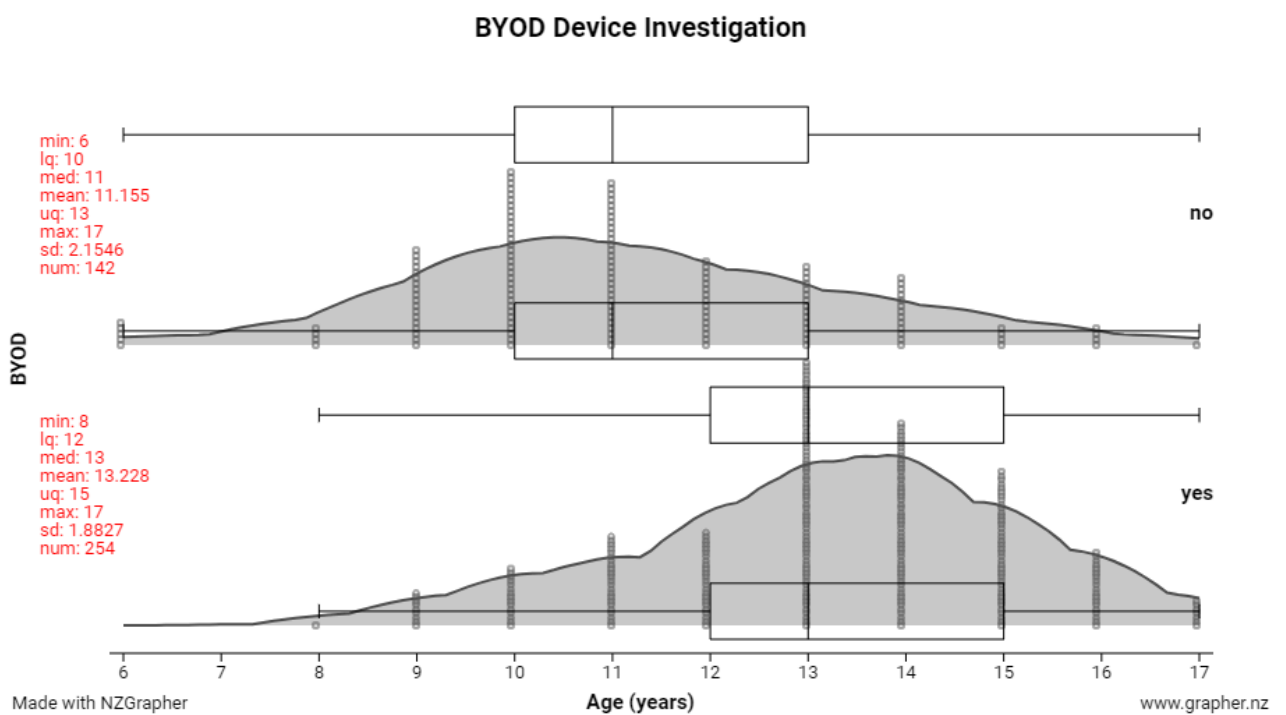


**For my sample**, I notice:

- The shape of the females and male bag weights have the same right skewed shape. The females and male bag weights are right skewed because they have one peak on the left-hand side, are asymmetric, and there is a longer tail on the right-hand side.
- The median of the female bag weights is a little heavier than the bag weights for males by 0.5kg. My evidence is that the median bag weight for females is around 3.8kg while the median bag weight for males is around 3.3kg.
- The spread of the middle 50% of females bag weights is slightly larger than the spread of males' bag weights, because the IQR of the females is approximately 3.5 kg compared to the IQR for males of 3.1 kg.




- 2) This investigation is looking into whether the ages of students who have a device is older than the age of students who don't have a device.





# Conclusion

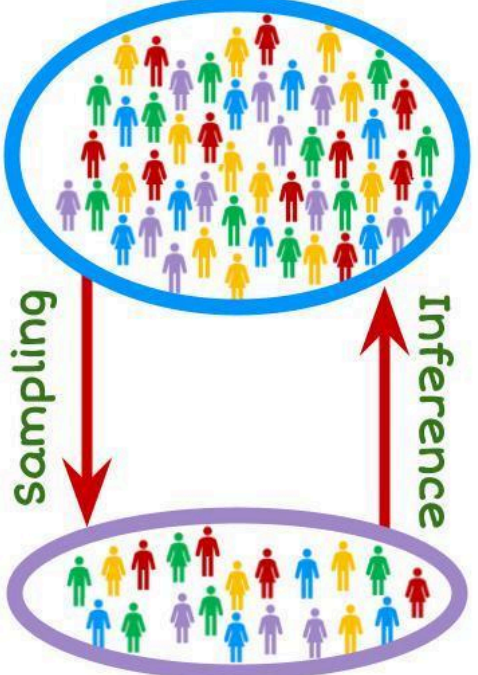


This section focuses on making inferences about the **Population**

© Liz Sneddon

## Inference

Drawing conclusions about the population based on a sample taken from the population.



© Liz Sneddon

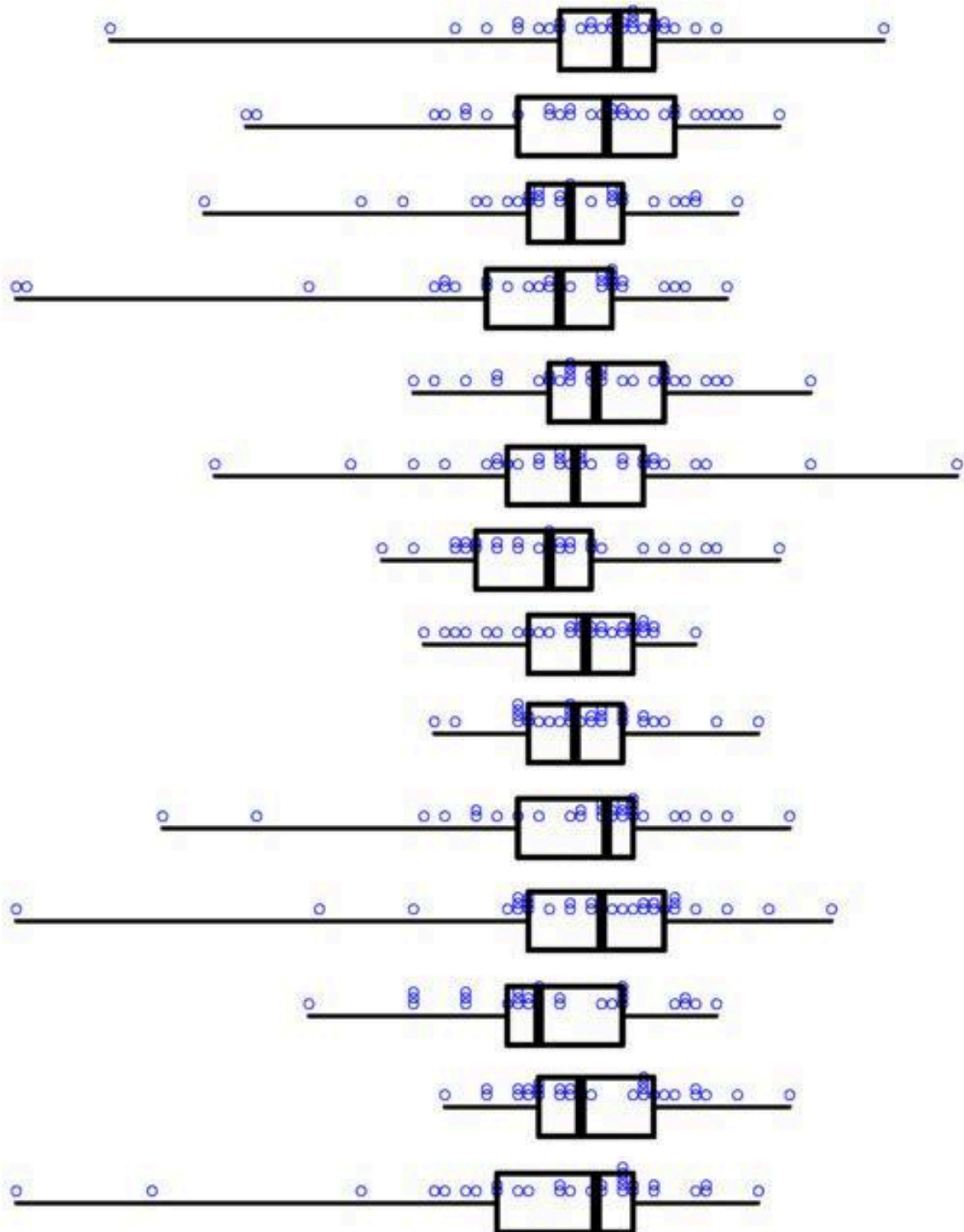
You need to include the following in your conclusion:

- Make an inference about the population using your confidence interval,
- Answer your investigation (can you make the call?),
- Discuss sampling variability.

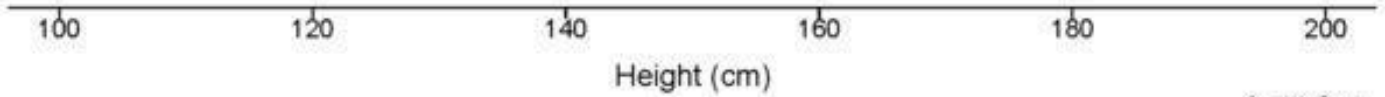
All of these things require an understanding of sampling variability, so we will start by exploring this in more depth.

## Exercise 16:

Below are graphs all of which are samples of size 30, from a population of height of 12-year-old girls in NZ.<sup>6</sup> Here is the animation: [bit.ly/VariationAnimation](http://bit.ly/VariationAnimation)



<sup>6</sup> Thanks to Prof. Chris Wild <https://www.stat.auckland.ac.nz/~wild/WPRH/>



*Class 10, June 09*

- a) Highlight the medians for each of the samples above.
- b) Each time we take a sample, would we randomly select the same people each time?

---



---



---

- c) Explain why the medians are different for each of the samples on the previous page.

---



---



---



---



---



---



---



---

- d) Find the smallest and largest estimate for the sample medians. Estimate how much do these values vary by, e.g., 1cm?

---



---



---



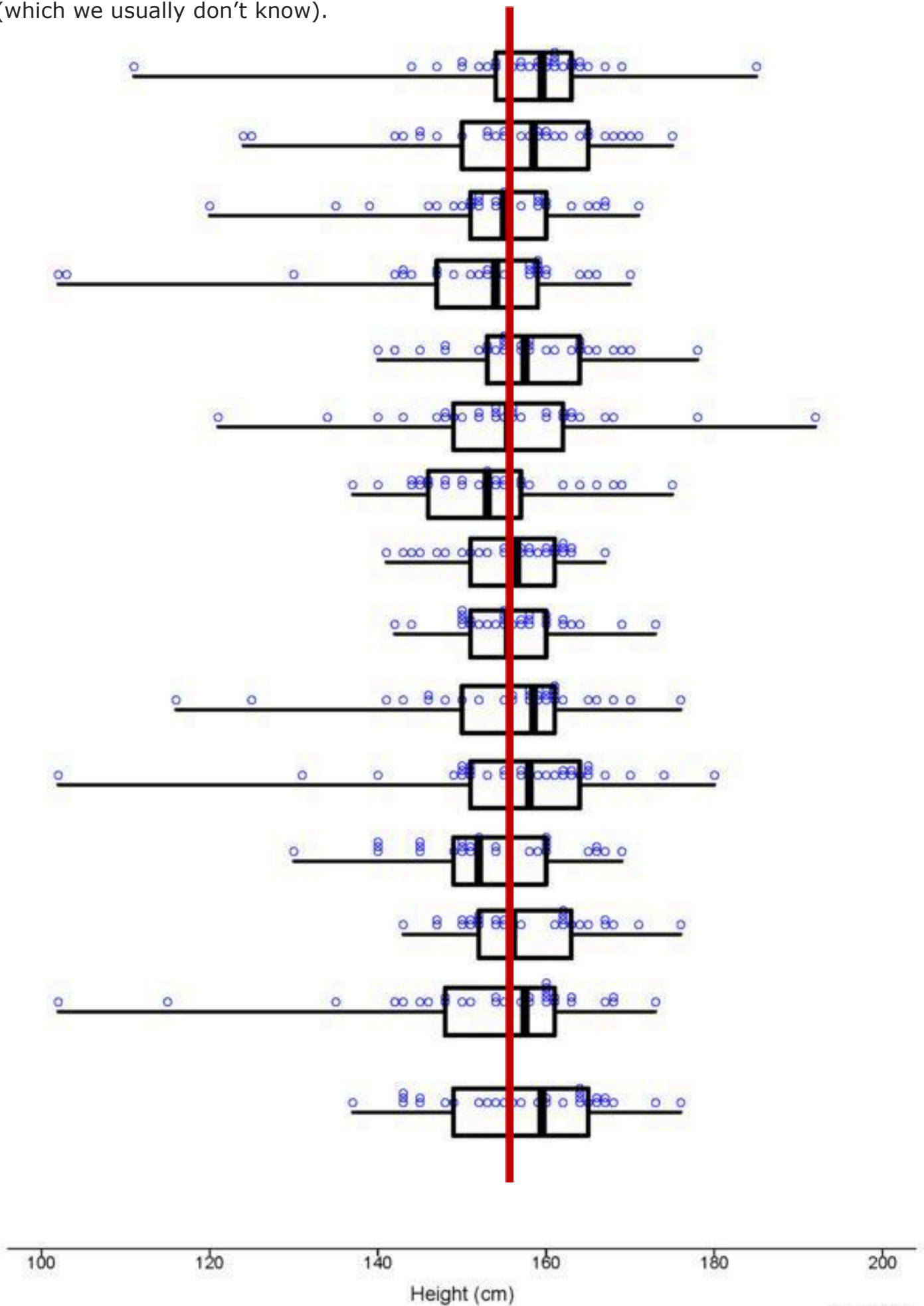
---



---



I've added the line drawn from the axis to the top, which is the population median (which we usually don't know).



e) What do you notice about the sample medians compared to the population medians?

---

---

---

---

---

---

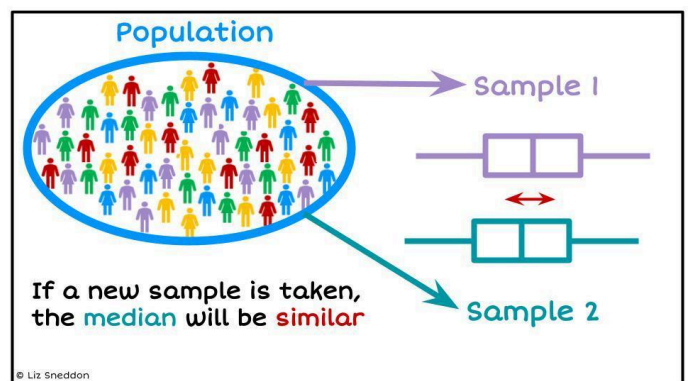
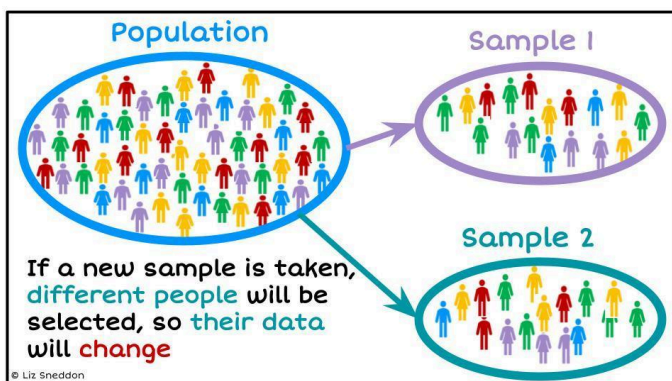
---

---

---

---

This leads us to a big idea, that the **sample median** is **likely** to be **close** to the **population median**, but our estimate is not exact.



# Sample size and sampling variation

## Exercise 17:

Go to the animation: [bit.ly/SamplingVariation](http://bit.ly/SamplingVariation)

This shows you different samples being taken, with different sample sizes. Look at the blue line which is the median. What do you notice about the amount of variation in the medians with the different sample sizes?

---

---

---

---

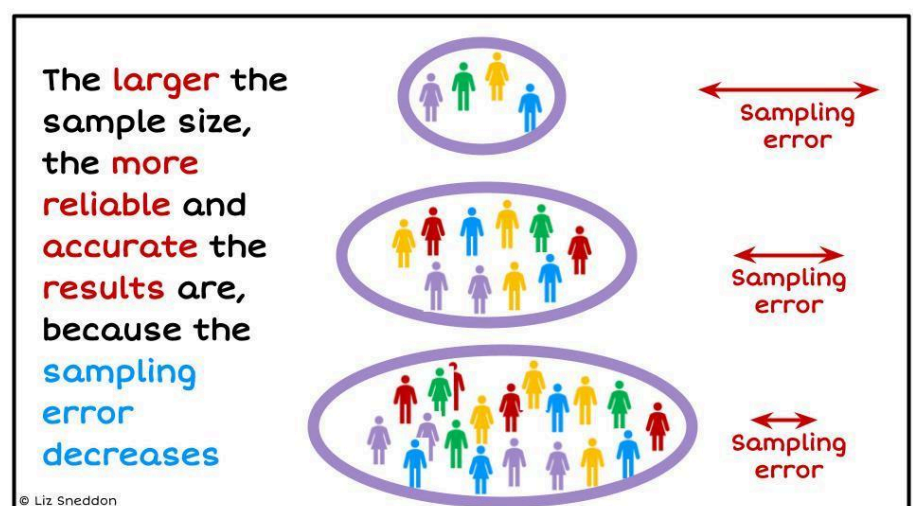
---

---

---

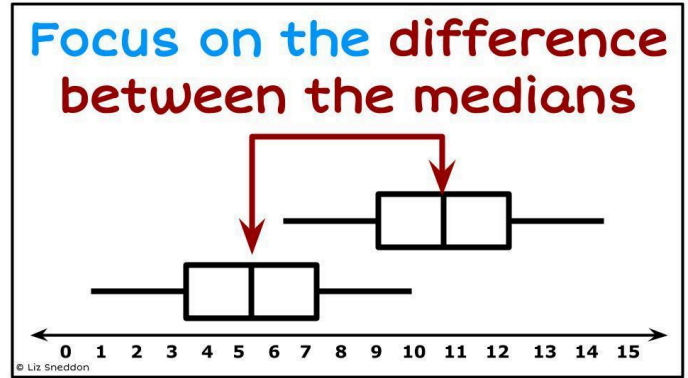
---

This leads us to a big idea, that the larger the sample size, smaller the **sampling error** is.



# Difference between the medians

The next step is to focus on the differences between the medians.

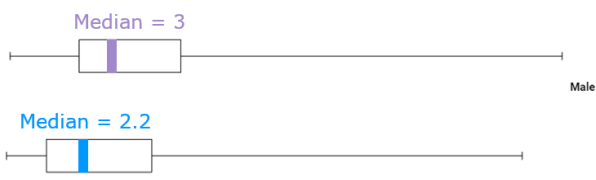
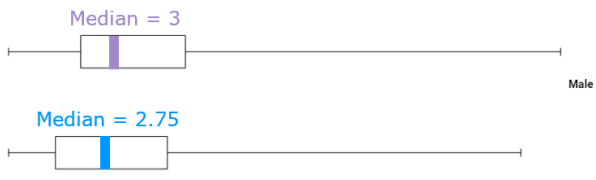
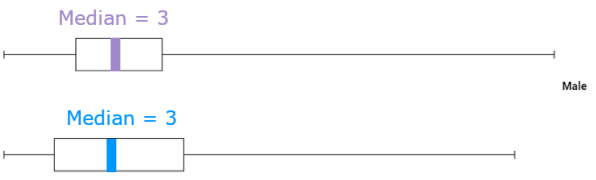
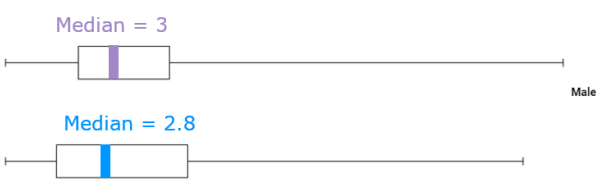


## Exercise 18:

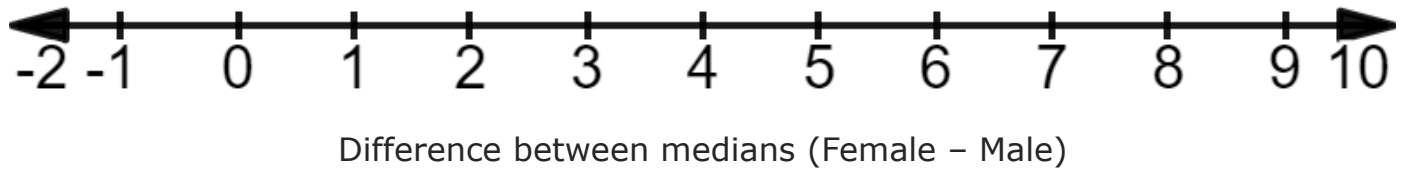
Each of the graphs below is a different sample of bag weights for teenagers at school in NZ. Calculate the difference between the medians for each graph.

**Difference = Female median – Male median**

1)	<p>Female Median = 3</p> <p>Male Median = 2.7</p>	2)	<p>Female Median = 3</p> <p>Male Median = 2.9</p>
Difference =		Difference =	
3)	<p>Female Median = 3</p> <p>Male Median = 2.95</p>	4)	<p>Female Median = 2.95</p> <p>Male Median = 2.3</p>
Difference =		Difference =	
5)	<p>Female Median = 3</p> <p>Male Median = 3</p>	6)	<p>Female Median = 2.85</p> <p>Male Median = 2.2</p>
Difference =		Difference =	

<p>7)</p>  <p>Difference =</p>	<p>8)</p>  <p>Difference =</p>
<p>9)</p>  <p>Difference =</p>	<p>10)</p>  <p>Difference =</p>

Next, draw add the differences that you calculated onto the axis below.



This leads us to understanding the basis of the **Bootstrap distribution**.

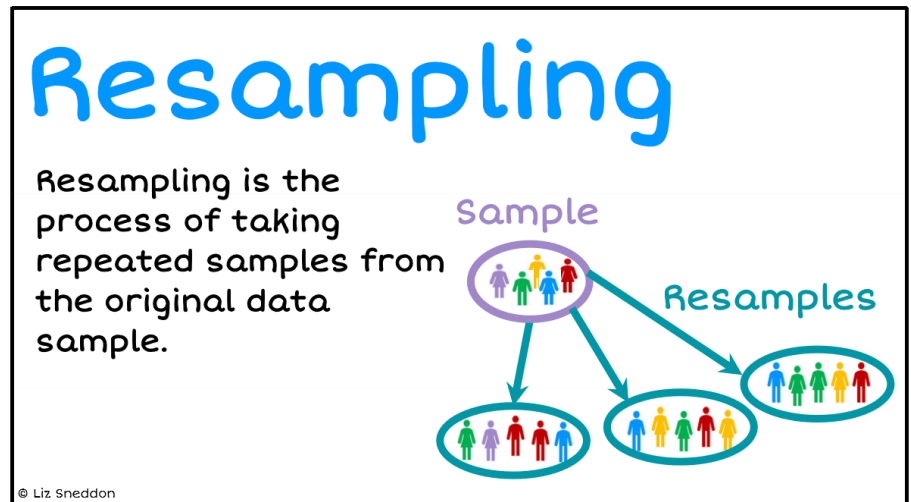
# Bootstrapping Distribution

The diagram below shows the process for a bootstrap distribution.

We start by having a population that we are interested in investigating, and we then take as big a sample as time and money allows.

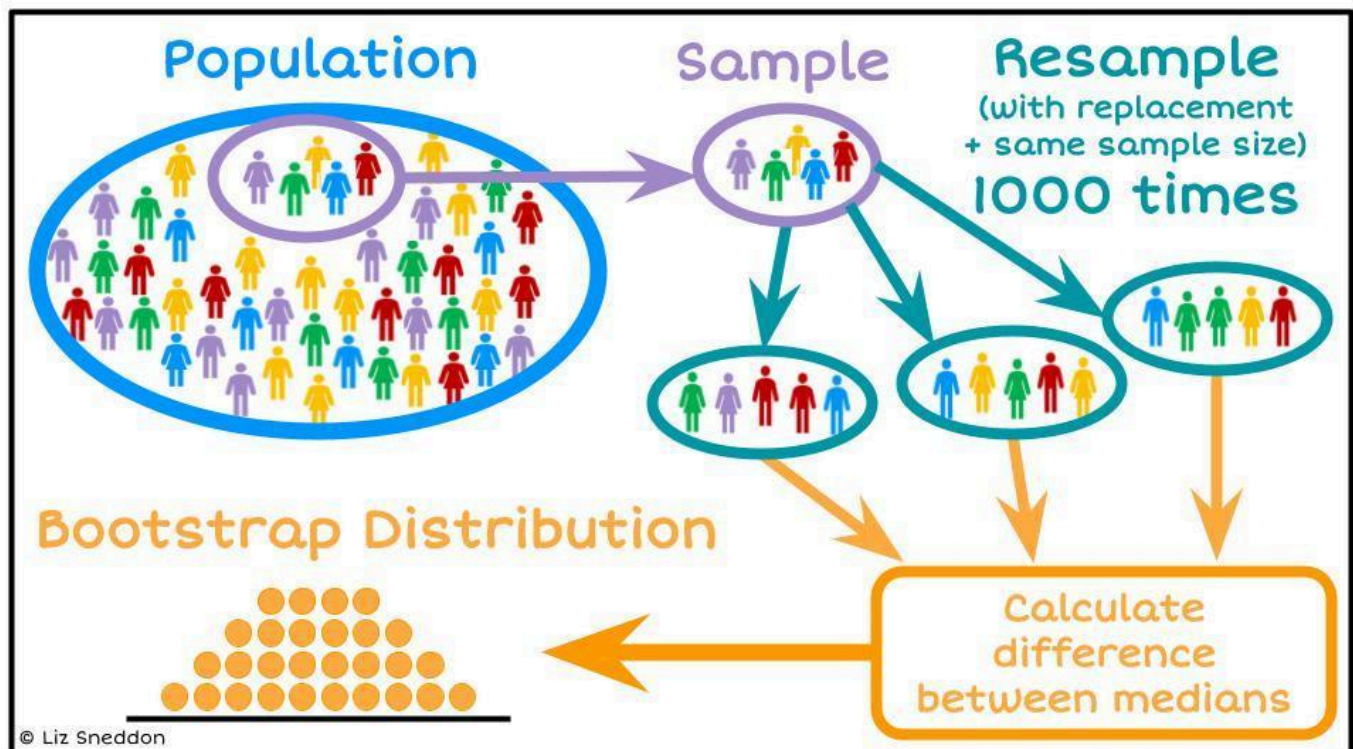
Then in order to improve the reliability and accuracy of the estimated results, we carry out a bootstrap process (this is a way of reusing the sample data to improve our estimates without needing to spend more money collecting samples - it does rely on the original sample being collected randomly and being representative of the population).

A re-sample is collected from the original sample, and people are selected randomly with replacement, where the sample size ( $n$ ) is the same as the original sample.

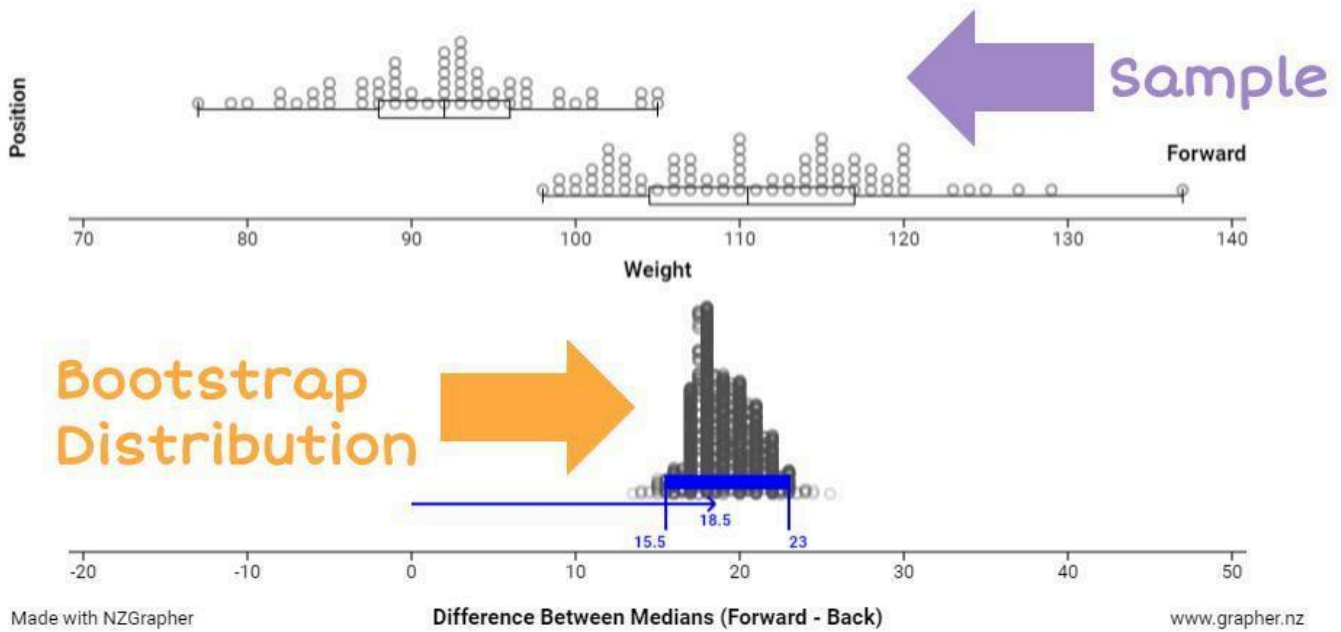


The difference between the medians of this re-sample is calculated and added to a dot plot.

This re-sampling process is repeated 1000 times to form the **Bootstrap Distribution**.

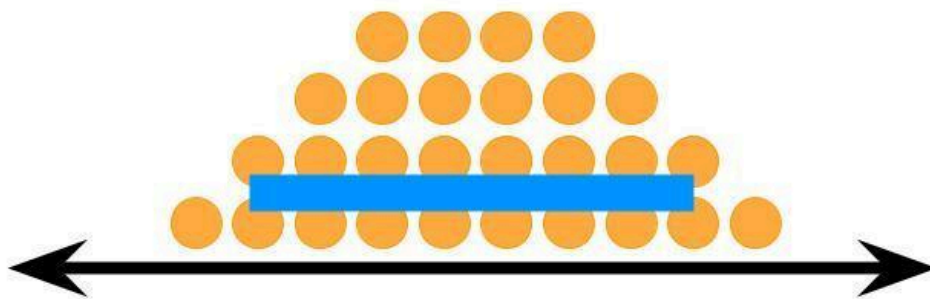


## Example:



## Bootstrapping Confidence Intervals

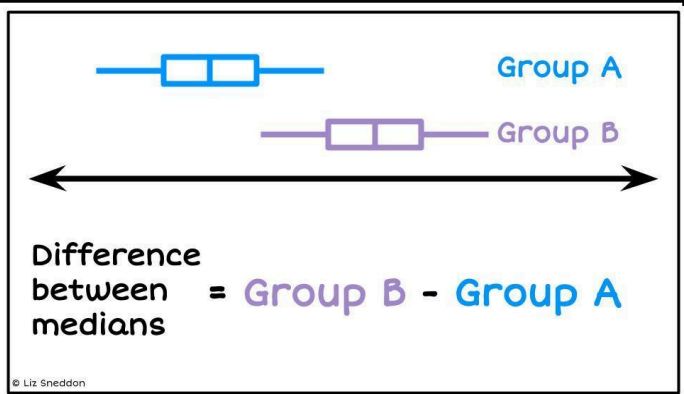
Once the Bootstrap distribution is formed, a confidence interval is then added.



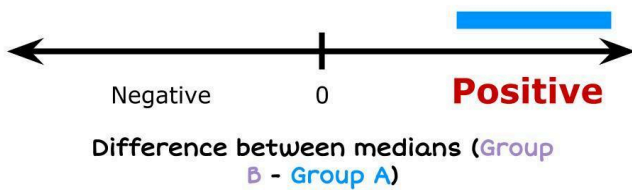
This is a confidence interval for the  
**difference between the medians**

A confidence interval is a range of values  
that **difference between the medians back**  
in the population is likely to fall within.

Always check the label on the axis underneath the bootstrap distribution. We will use this to work out what a positive or negative confidence interval limit means.



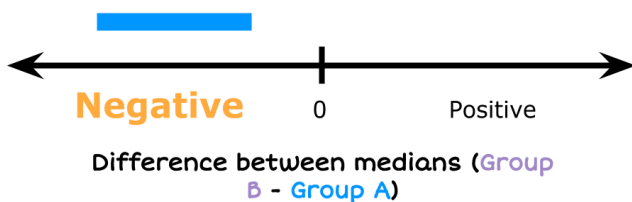
If the bootstrap confidence interval for the difference is **positive**, Group B tends to be **larger** than Group A



E.g.

Group B median is likely to be between \_\_\_\_\_ and \_\_\_\_\_ **more** than Group A median, for the population.

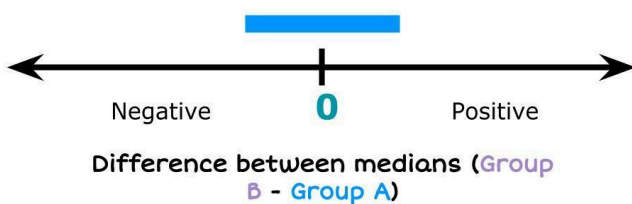
If the bootstrap confidence interval for the difference is **negative**, Group B tends to be **smaller** than Group A



E.g.

Group B median is likely to be between \_\_\_\_\_ and \_\_\_\_\_ **less** than Group A median, for the population.

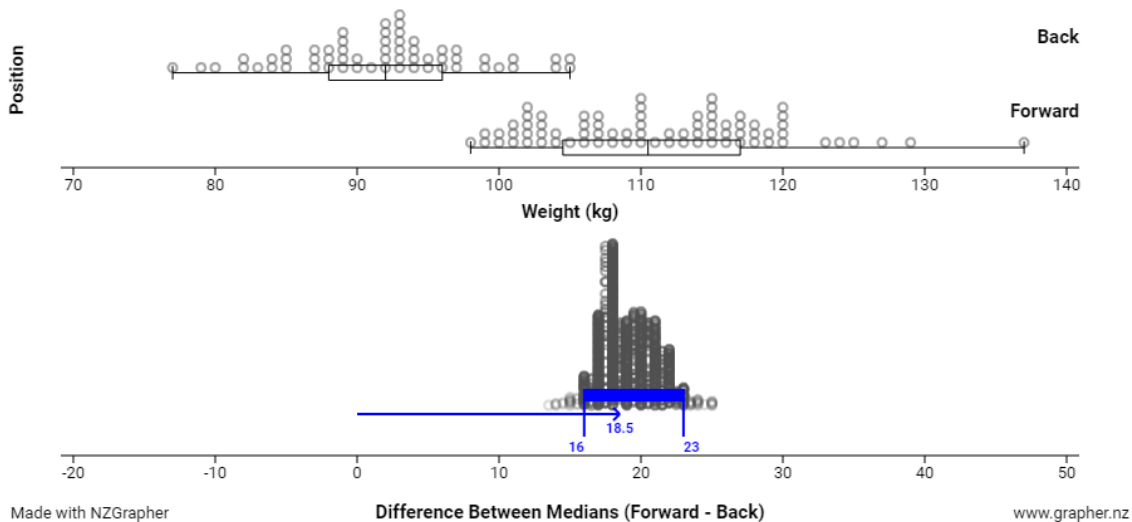
If the bootstrap confidence interval for the difference **includes zero**, Group B could be **larger or smaller** than Group A



E.g.

Group B median is likely to be between \_\_\_\_\_ **less** and \_\_\_\_\_ **more** than Group A median, for the population.

## Example 1:

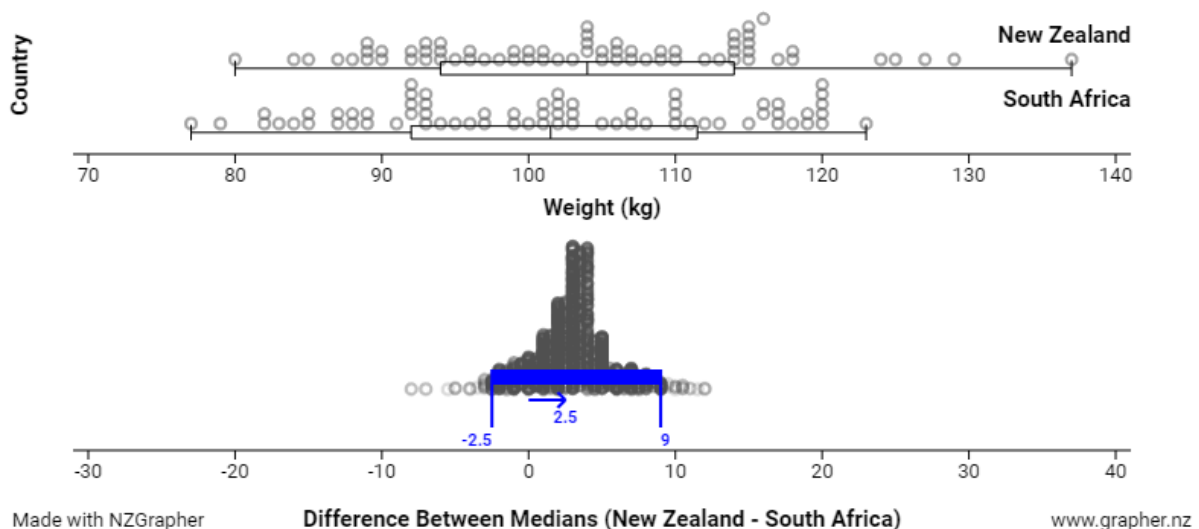


From the bootstrapping confidence interval, I am fairly sure that the median weight of forwards will be between 16 kg and 23 kg more than median weight of backs for the **POPULATION** of ALL top players listed on the website <http://www.rugby-sidestep-central.com/>.

## Example 2:

### Problem:

What is the difference between the median weight (kg) of rugby players in New Zealand and South Africa from ALL top players listed on the website <http://www.rugby-sidestep-central.com/>



From the bootstrapping confidence interval I am fairly sure that the median weight of New Zealand rugby players could be up to 9 kg more or up to 2.5kg less than median weight of South African rugby players, for the **POPULATION** of ALL top players listed on the website <http://www.rugby-sidestep-central.com/>









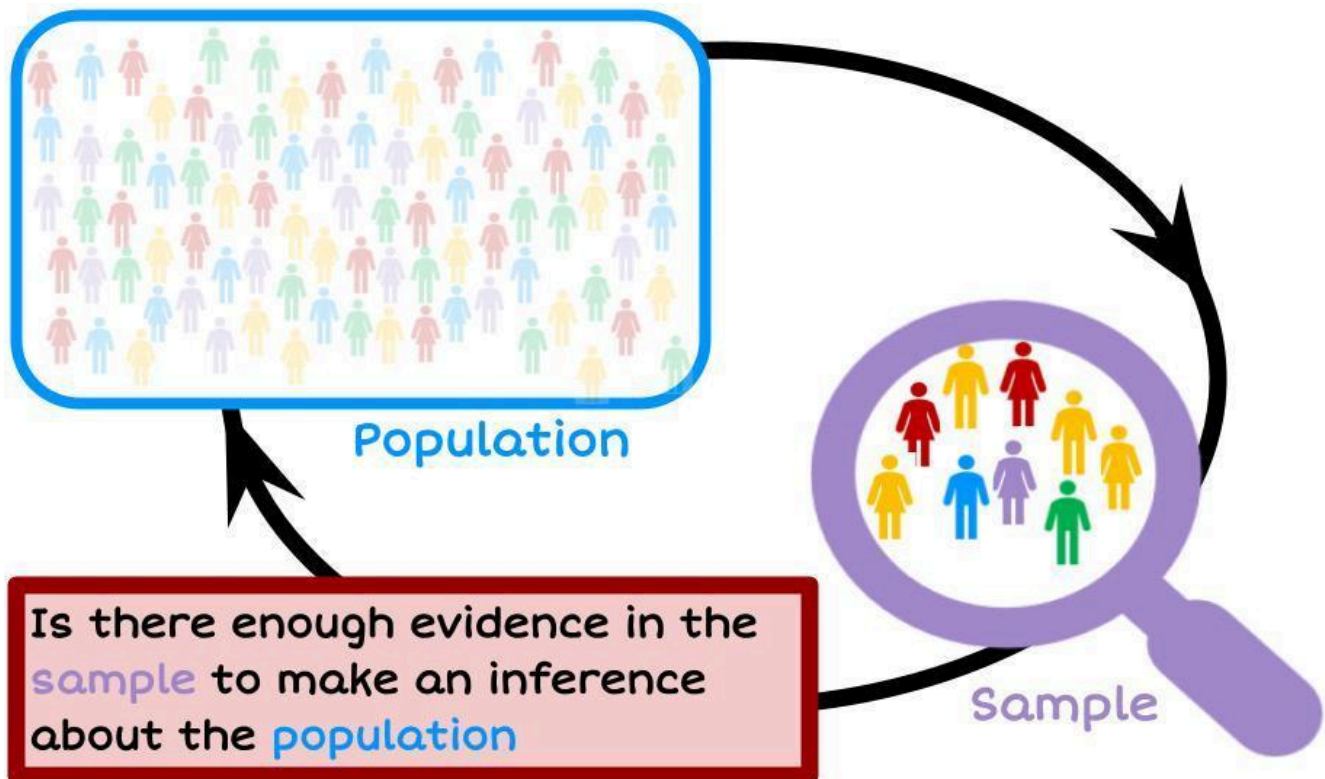






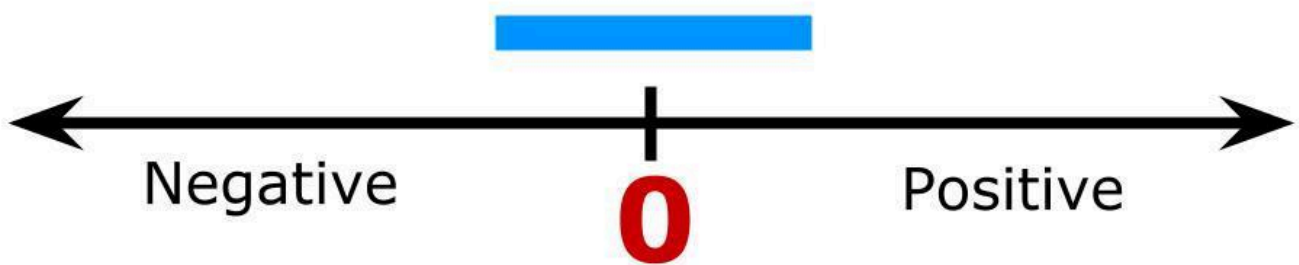
# Is there enough evidence?

We are looking to see if we have **enough evidence** of a difference between the medians, **back in the population.**



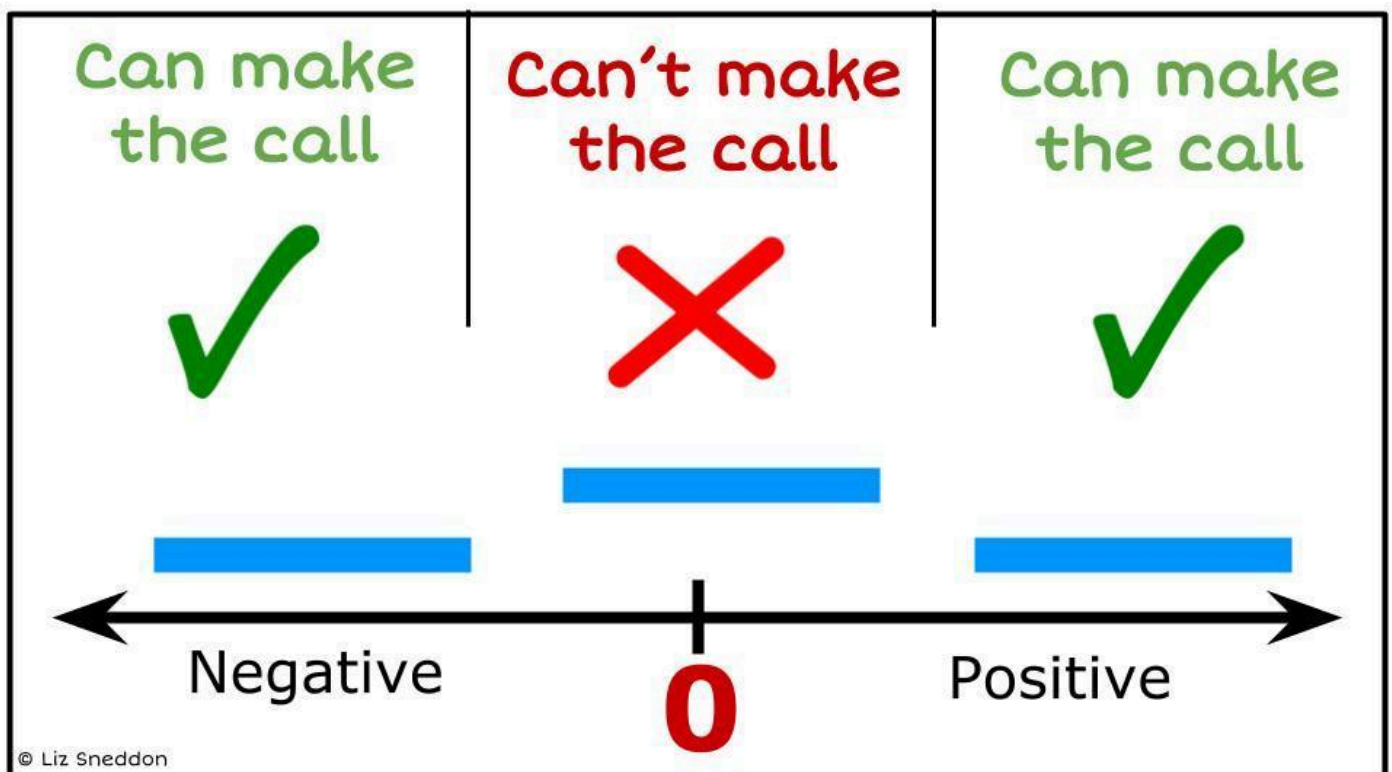
## Making the call

To decide if you can make the call or not, look at whether the bootstrap confidence interval **CONTAINS ZERO**



© Liz Sneddon

If the bootstrap confidence interval **contains zero**, this suggests that the **medians could be the same**, therefore we **can't make the call**.



© Liz Sneddon

# Answering the investigation question

## To answer the comparison question you need:

- Categorical variable,
- Numerical variable,
- Compare the medians,
- Direction,
- Population **(ALL)**.

© Liz Sneddon

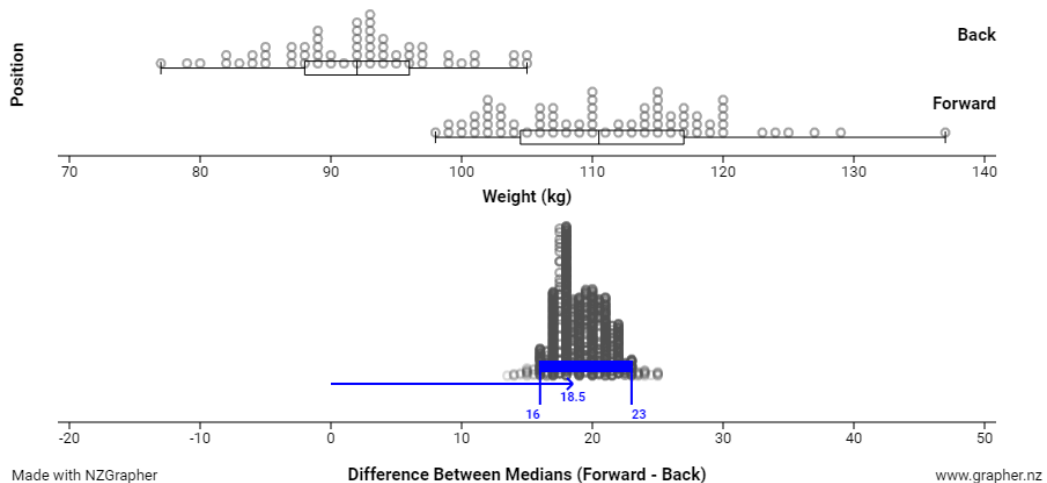
Don't forget to include the **direction** in your answer.

### Example:

#### Problem:

What is the difference between the median weight (kg) of forward and back rugby players from ALL top players listed on the website

<http://www.rugby-sidestep-central.com>.



#### Conclusion:

The bootstrap confidence interval **DOES NOT** contain zero, so we **CAN** make the call.

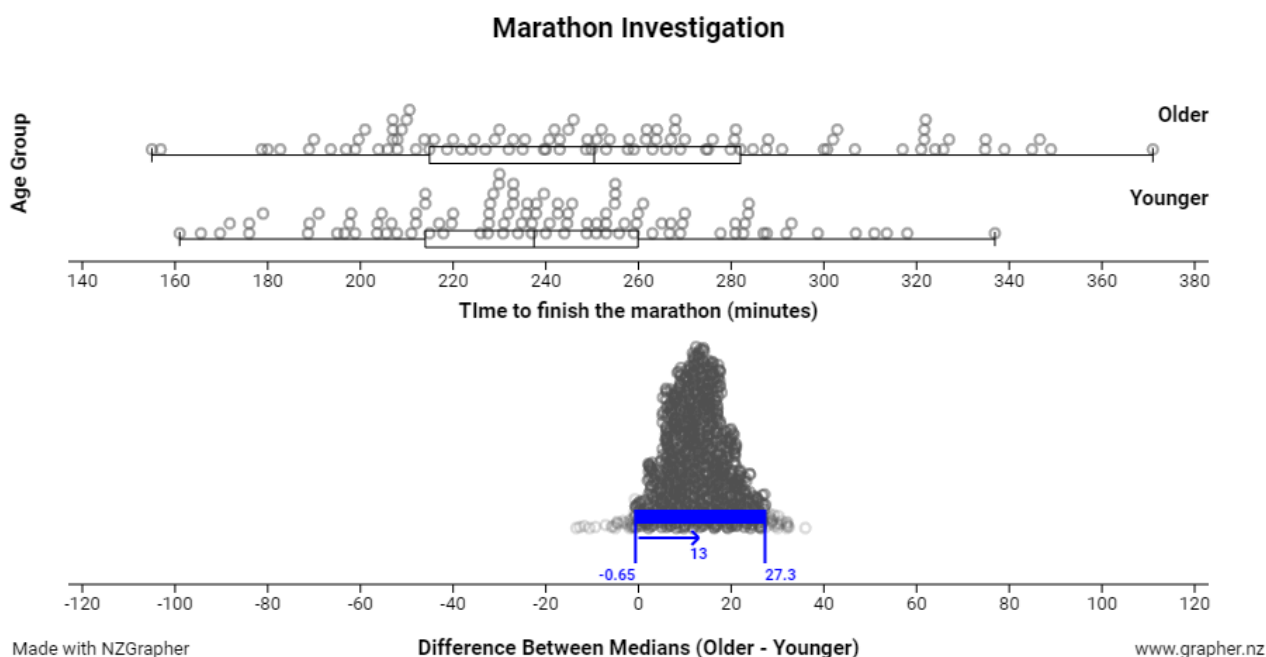
I can make the call, so I **DO** have enough evidence that the median weight of forwards is heavier than the median weight of backs, for **ALL** top players listed on the website <http://www.rugby-sidestep-central.com>.

## Exercise 20:

For the investigations below, answer the investigation question, justifying your conclusion.

### 1) **Problem:**

What is the difference between the median length of time it takes older (over 40 years old) and younger (under 40 years old) people to complete a marathon, for ALL marathon runners in NZ?



Does the CI include zero?	Yes / No
Can you make the call?	Yes / No

Answer the investigation question:

---



---



---



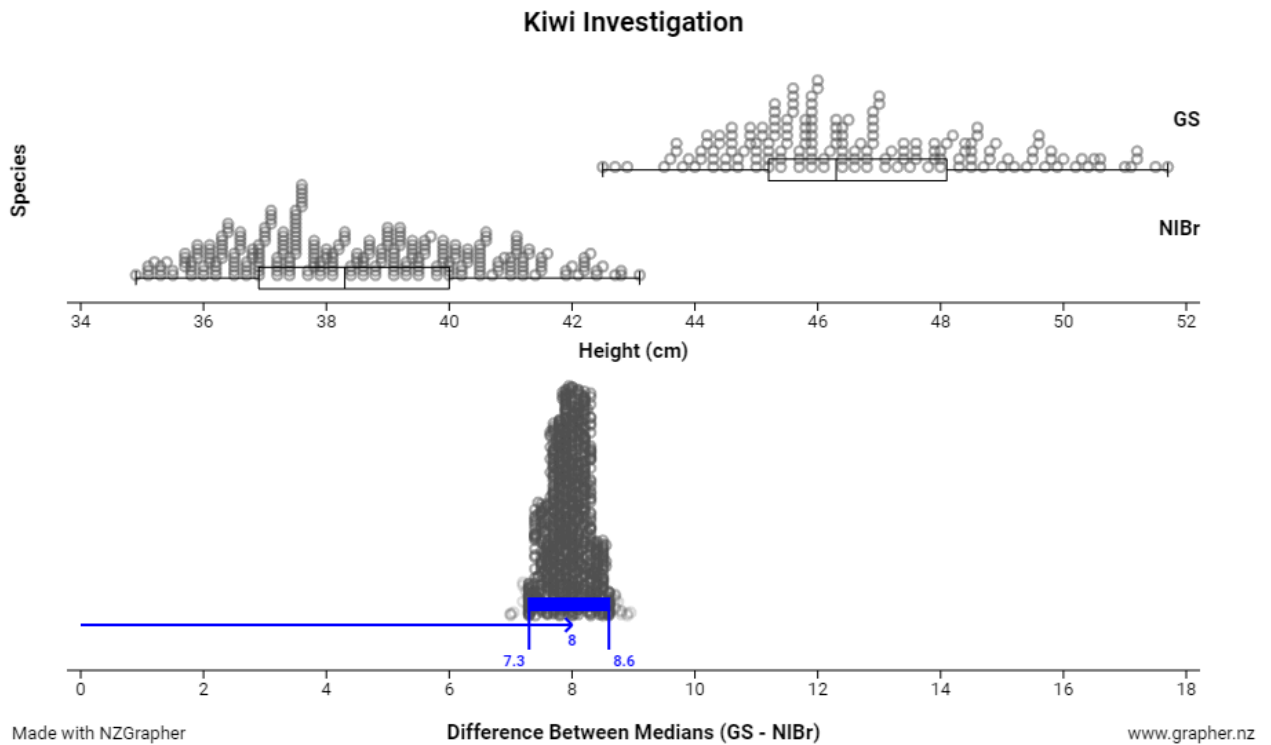
---



---

2) **Problem:**

What is the difference between the median height of the Great Spotted kiwi and the North Island Brown kiwis, for all kiwi birds in NZ?



Does the CI include zero?	Yes / No
Can you make the call?	Yes / No

Answer the investigation question:

---

---

---

---

---

---

---

---

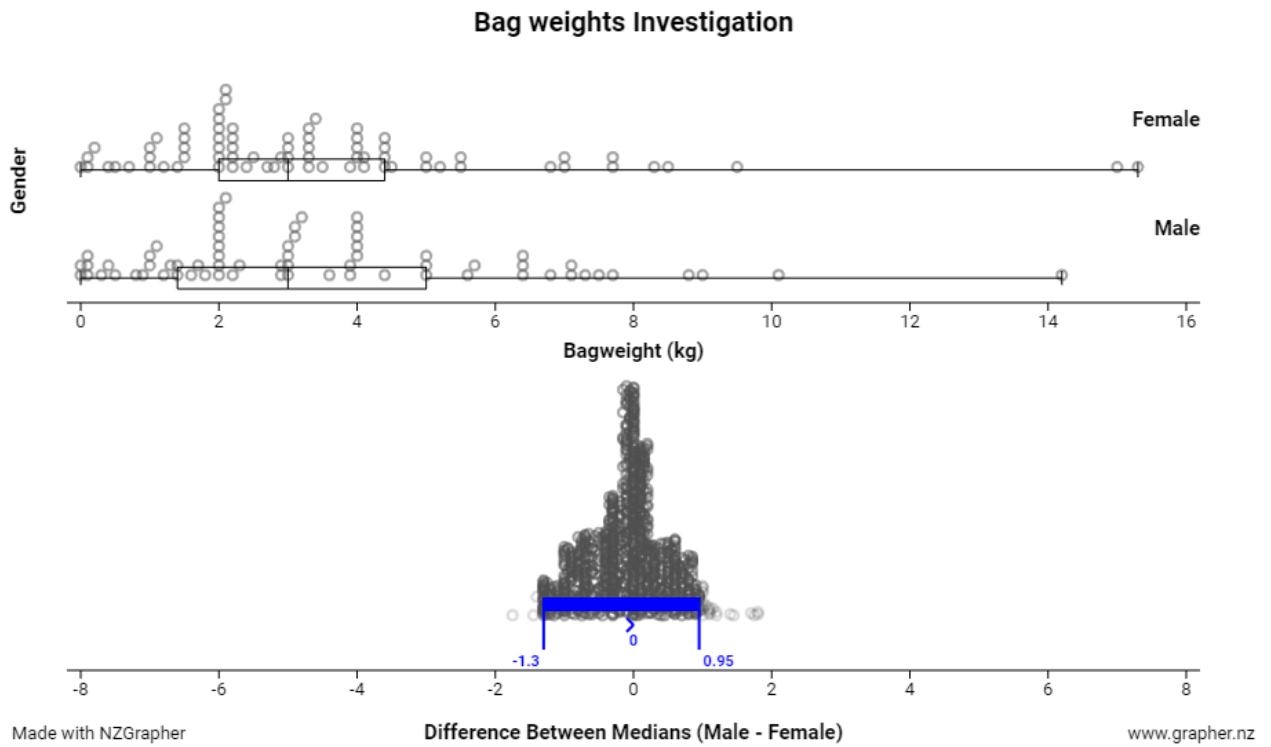






6) **Problem:**

What is the difference between the median weight of school bags for female and male teenagers in NZ?



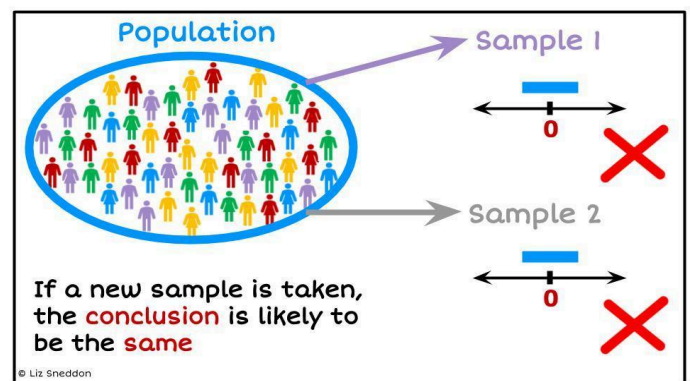
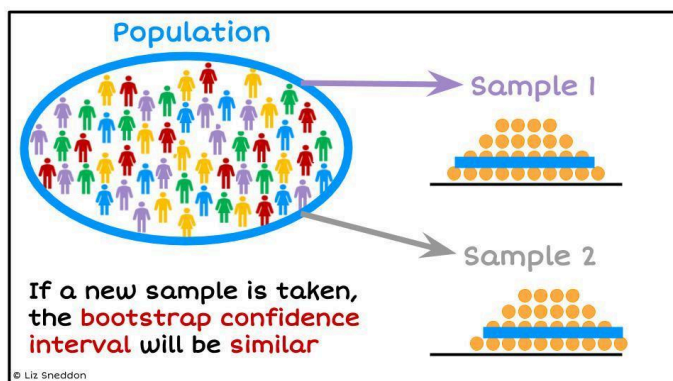
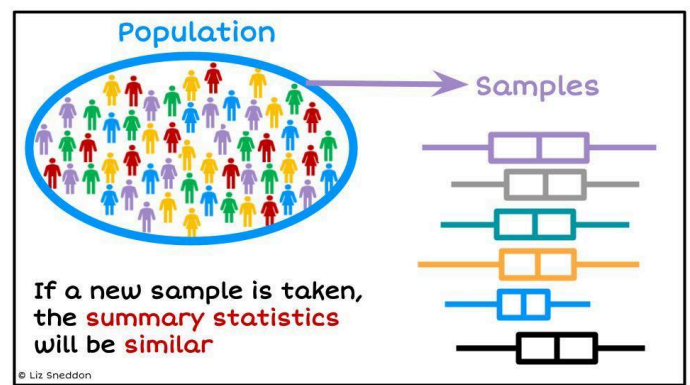
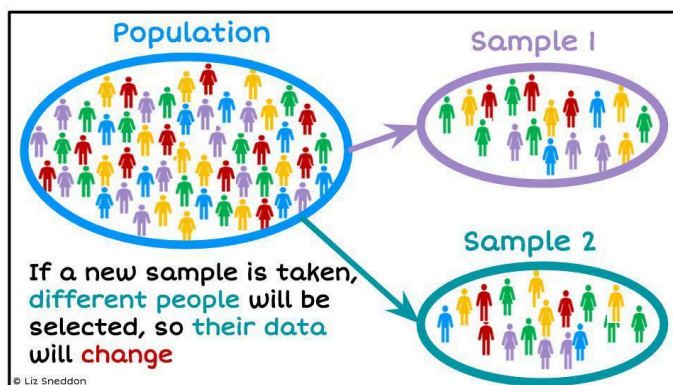
Answer the investigation question:

# Sampling variability

This describes the variation (differences) that happens when you take another sample. You always have to write an explanation in context.

You need to consider:

- how would the **data** change?
- how would the **analysis** change (particularly the summary statistics)?
- how would your **bootstrap confidence interval** change?
- how would your **conclusion** change?



## Example:

If I took another sample, I would get different weights for rugby players as I would be collecting data from different rugby players. I would expect though that the summary statistics (minimum, LQ, median, UQ and maximum) weights of forwards and backs to be similar to the values in my sample. Because the median weights for forwards

and backs would be similar, this would lead to a similar bootstrap confidence interval, and therefore the conclusion that the weights of forwards are larger than the weights of backs, is likely to stay the same.



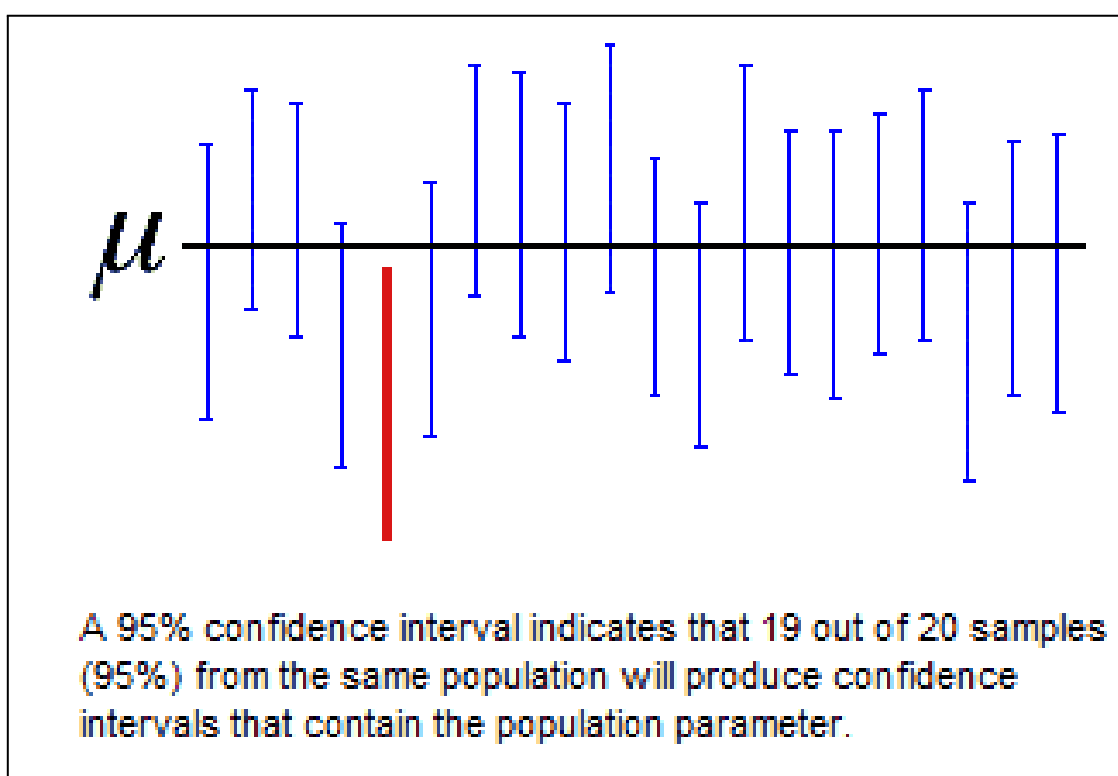


# Confidence intervals and population medians

---

The second idea we need to cover is whether or not in **our one sample**, how likely we are to have captured the true difference between the population medians in our confidence interval?

The image below shows twenty different confidence intervals, and out of these twenty intervals, there is only one that does not cover the population median.



When we take a sample, we don't know which of these confidence intervals we have, and whether or not we have the one that does not contain the difference between the population medians. However, because most of the confidence intervals **DO** contain the difference between the population medians, it is quite likely that our confidence interval **does** contain the difference between the population medians.

# Other Factors

---

Think about making a puzzle. If you have one piece of the puzzle, you can learn a lot - the colours, the shapes, and guess where it might go in the big picture. But you know that this is only one piece of the puzzle, and that there are lots of other pieces that when put together make a whole different picture.

Data is like a puzzle. You might investigate one particular piece, one particular comparison, but we always need to be aware that there are many other pieces in this puzzle. Before we can paint the whole picture, we need to think about what the other pieces might be.

This is the idea of other factors. You need to consider what other factors might also affect the numerical variable that you are interested in and use research to explain **WHY** this factor might be affecting it.



## Example:

---

### **Problem:**

What is the difference between the median weight of teenagers than young children, for all children in NZ.

### **Other variables or factors:**

There are a number of different factors that might affect the weight of children in NZ. For example, if a child has parents who are both slim and short in stature, then because of the genetic link it is likely that the child is also likely to be slim and short in stature. Equally, a child whose parents have bigger and heavier bones, and a wide/tall build are likely to be taller and heavier. A study published in the UK<sup>7</sup> supports this and discusses how there is a link between people's weight and their genetics, where being slim is a heritable trait.

Another factor that could affect the weight of children is the amount of exercise they do each week. I expect that a child who is more active and spends more time each week exercising would have less body fat than a child who is less active and spends less time exercising each week. The more body fat a child has the higher their weight will be. Kids Health<sup>8</sup> suggest that "kids can reach a healthy weight by eating right and being active".

---

<sup>7</sup> <https://www.healthline.com/health-news/heres-how-much-your-genes-impact-your-ability-to-lose-weight>

<sup>8</sup> <https://kidshealth.org/en/parents/childs-weight.html>





