

Reviewing Statistics

Throughout your study of strand 1 you will have considered all aspects of a statistical approach:

- asking a question that results in data that varies
- displaying this data in a way that allows you to see patterns in the variation
- analysing the patterns in the data
- drawing conclusions from that data.



You may even have had an opportunity to get a glimpse of what it is like to become a statistical detective; attempting to account for unexpected variability you observe in a particular set of data.

As you review for the final examination in June, it is important that you can connect each element of your study and consider the BIG IDEA of the strand so that you will be able to use the elements appropriately to help you solve problems that you may not have seen before.

The following is an extract from Strand 1 of the syllabus; it summarises what you should be able to do when you finish studying this strand.

It is envisaged that throughout the statistics course students will be involved in identifying problems that can be explored by the use of appropriate data, designing investigations, collecting data, exploring and using patterns and relationships in data, solving problems, and communicating findings. This strand also involves interpreting statistical information, evaluating data-based arguments, and dealing with uncertainty and variation.

You may decide to form a study group with your friends or you may prefer to work alone; either way as you work through this review document you will consider issues such as framing a question in order to obtain meaningful **reliable** data, selecting a sample in order to avoid **bias**, **displaying** your data in a way that will allow you to see patterns in the variation and **drawing conclusions** from your data.

Asking the Question

Think



Do you use a computer?

How did you answer the question?

What were you thinking when you answered it?

A university sports outlet was considering shutting down their campus shop and becoming an on-line store in an effort to reduce costs. A group of students was surveyed and asked that same question:

Do you use a computer?

Sophie answered **Yes** because she thought the question meant had she ever used a computer.

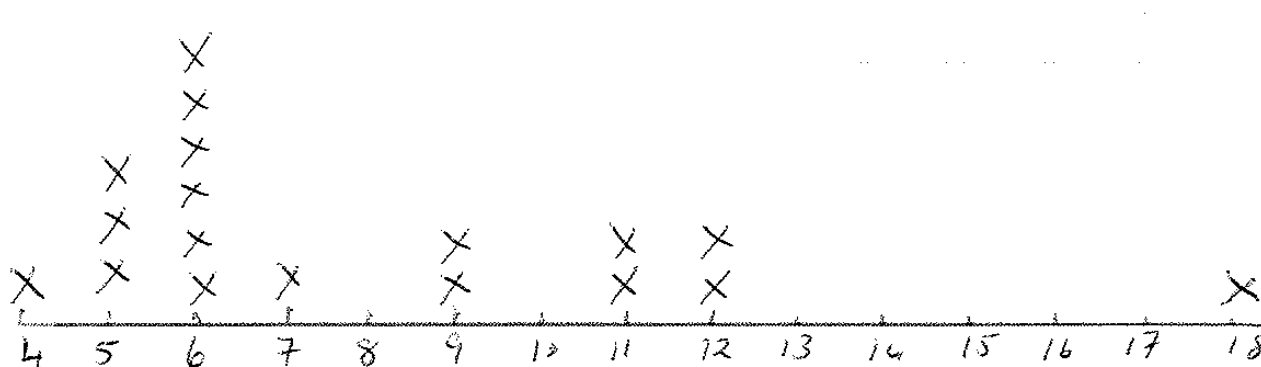
Joe answered **No** because he thought the question was asking whether he used one regularly.

Andrew answered **No** because he played games on the computer and didn't think this counted as "using" one.

Do you think the results of this survey are **reliable**?

How could you rephrase the question so that it is less ambiguous and more likely to provide useful answers?

A group of students interested in finding the typical family size for their class obtained the data displayed in this line plot



What question do you think they asked in order to elicit this data?

What issues would they have needed to consider when framing the question?

Displaying the data and drawing conclusions from it

Use fractions or percentages to describe the data.

Can you see any clumps or areas where a large proportion of the data falls?

Are there any unusual family sizes? [18 is an unusual value in this set.]

What do you think is the typical family size of this group? Why?

If you were asked to predict the family size of someone from this group what value would you give?

Why?

How certain would you be? Can you lower this to a smaller range? How **confident** would you be now?

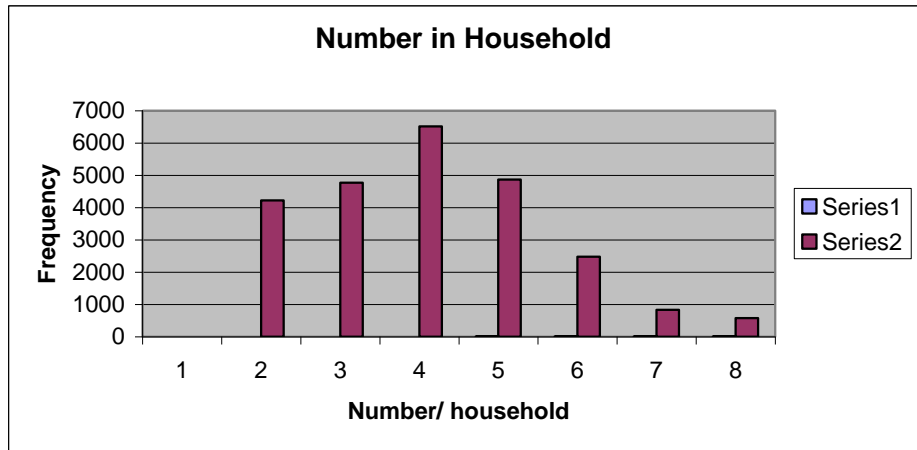
Calculate the mean family size for this group and identify the median family size. Which is a more reasonable estimate of typicality?

You could do a similar survey of your class, display the data in a line plot and compare the two data sets.

Or you could visit

<http://beyond2020.cso.ie/Census/TableViewer/tableView.aspx?ReportId=109241> and retrieve some data from the area in which you live, use Excel to display the data and compare it to the sample above.

Household size	Frequency
2	4218
3	4773
4	6512
5	4870
6	2478
7	833
8	572

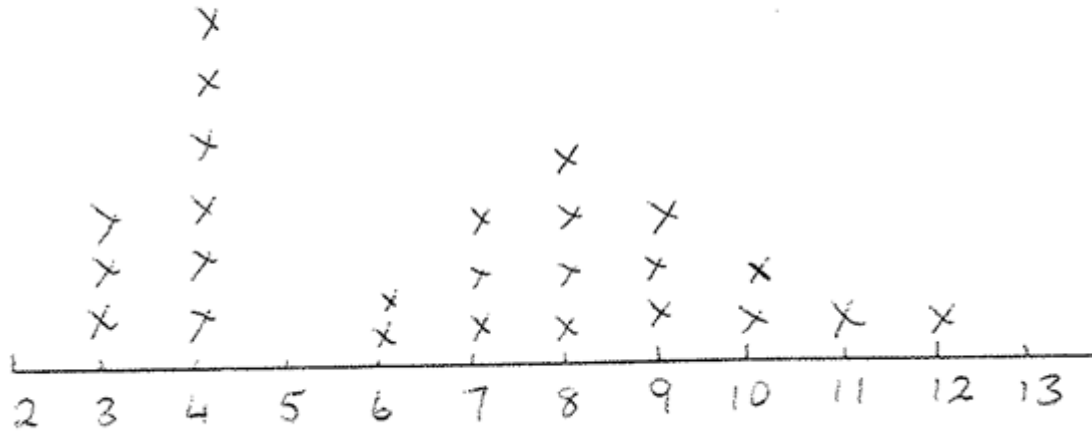


This bar chart was drawn with data from Carlow.

Compare this data with that from the sample set above. What is the range of this data set? What is the range from the sample data set?

Is there any evidence to suggest that the sample was from Carlow? Explain.

The following data set was gathered from a TY class who were interested in finding out what was the typical amount of money spent by their parents on the lotto each week.



Use fractions or percentages to describe the data.

Can you see any clumps or areas where a large proportion of the data falls?

Are there any unusual amounts?

What do you think is the typical amount spent on the Lotto each week by this group ? Why?

If you were asked to predict the amount spent on the Lotto each week by someone from this group what value would you give? Why?

How certain would you be? Can you lower this to a smaller range? How **confident** would you be now?

Return to the value you think is the typical amount spent on the Lotto each week by this group.

Calculate the mean amount spent on the Lotto by this group and identify the median amount spent on the Lotto each week. Which is a more reasonable estimate of typicality?

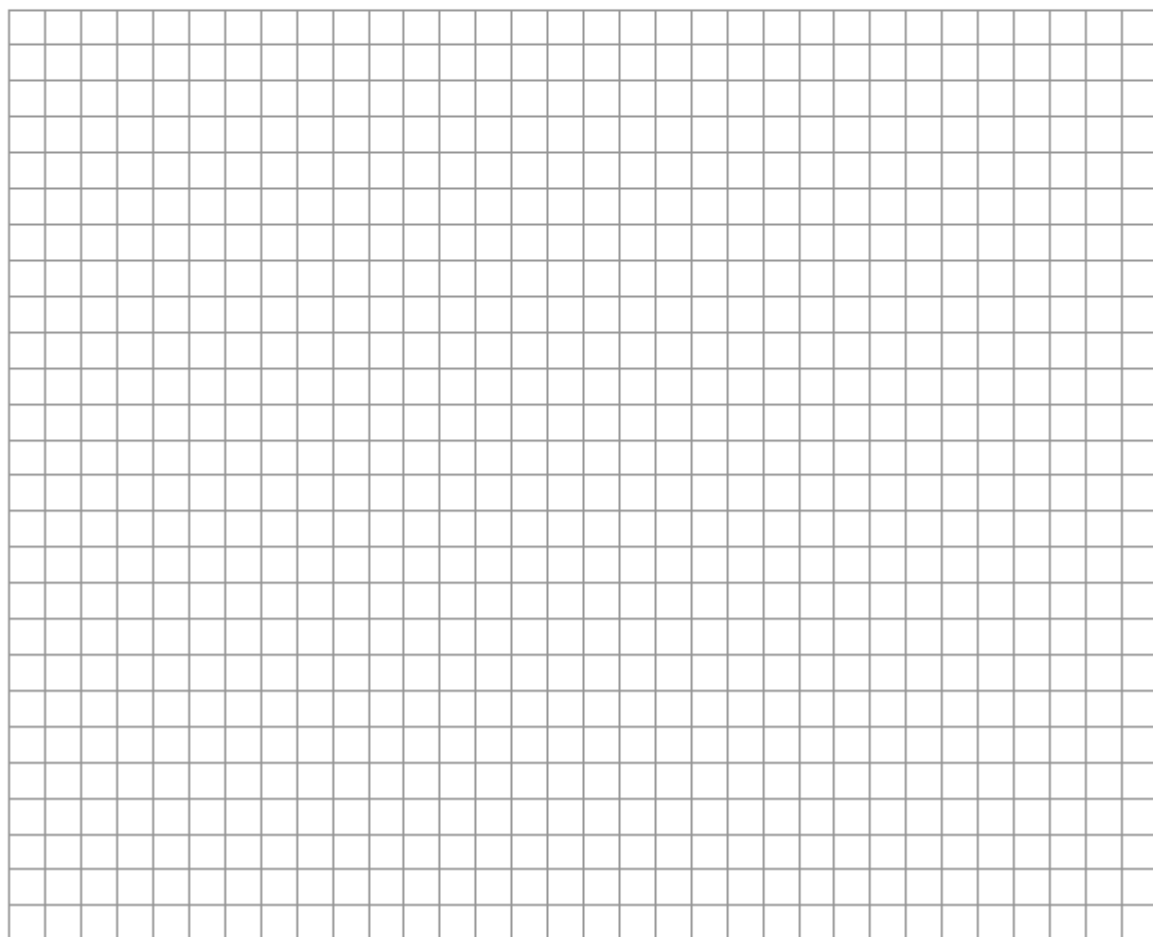
Q7 a) and b) on the **2010 HL** mock paper

Some research was carried out into the participation of girls and boys in sport. The researchers selected a simple random sample of fifty male and fifty female teenagers enrolled in GAA clubs in the greater Cork area. They asked the teenagers the question: *How many sports do you play?*

The data collected were as follows:

Boys	Girls
0, 4, 5, 1, 4, 1, 3, 3, 3, 1,	3, 3, 3, 1, 1, 3, 3, 1, 3, 3,
1, 2, 2, 2, 5, 3, 3, 4, 1, 2,	2, 2, 4, 4, 4, 5, 5, 2, 2, 3,
2, 2, 2, 3, 3, 3, 4, 5, 1, 1,	3, 3, 4, 1, 6, 2, 3, 3, 3, 4,
1, 1, 1, 2, 2, 2, 2, 2, 3, 3,	4, 5, 3, 4, 3, 3, 3, 4, 4, 3,
3, 3, 3, 3, 3, 3, 3, 3, 3, 3	1, 1, 3, 2, 1, 3, 1, 3, 1, 3

(a) Display the data in a way that gives a picture of each distribution.





Stop and think

Under what conditions would a **line plot** be a meaningful representation?

Under what conditions would a **stem and leaf plot** be a more meaningful representation?

Try using statistics to solve this problem.

PROBLEM: *Climbing helmets are made in a variety of styles and sizes.*

The manager of You Climb Safely must decide what styles of helmet to keep in stock and how many helmets of each size to order. A standard fit helmet is offered in 10 sizes. When you order helmets you must order 1000. How many of each helmet size should the manager order?

In order to get an idea of how head sizes are **distributed** the manager decided to measure the head circumferences of a group of people.

Think: what is the **population** of interest? Can he measure the circumferences of the heads of the whole population? How will he choose a **sample**?

The manager chose a **Simple Random Sample** of climbers from clubs around the country and recorded their head circumference and gender in the table overleaf.

Is this a suitable sample? Why or Why not?

Gender	Head Circumference (mm)
F	522
M	580
M	552
F	531
M	563
F	546
F	545
M	545
M	545
M	568
F	560
M	613
F	555
F	573
M	585
F	584
M	600
M	595
M	593
F	590
M	594
F	564
F	536
M	586
F	540
M	585
M	550
M	565
F	600
F	590
F	551
M	590
M	580
F	577

Is a line plot a good representation of this data?

Display the data in a stem and leaf plot.

Describe the data.....Are there any clumps or areas where the data is concentrated? Are some head sizes more common than others?

Use your representation to answer the original question: **how many helmets of each size should the manager order?**

Begin by counting the number of leaves on each stem.

Look at the first stem...52 ..How many leaves are there on stem 52? What fraction of the total is this? What % of the total number of head circumference measurements does stem 52 represent?

How many helmets size 520cm- 530cm should the manager order?

Continue working like this until you have decided how many helmets of each size the manager should order.

Return to your representation...Do you think there are **gender effects?** Try representing the male and female data in **back to back stem plots** Compare the two sets of data; is there any evidence to suggest that there are differences in the sizes of heads of men and women?

If there are gender effects will this affect the number of helmets the manager should order? Or are helmets unisex?

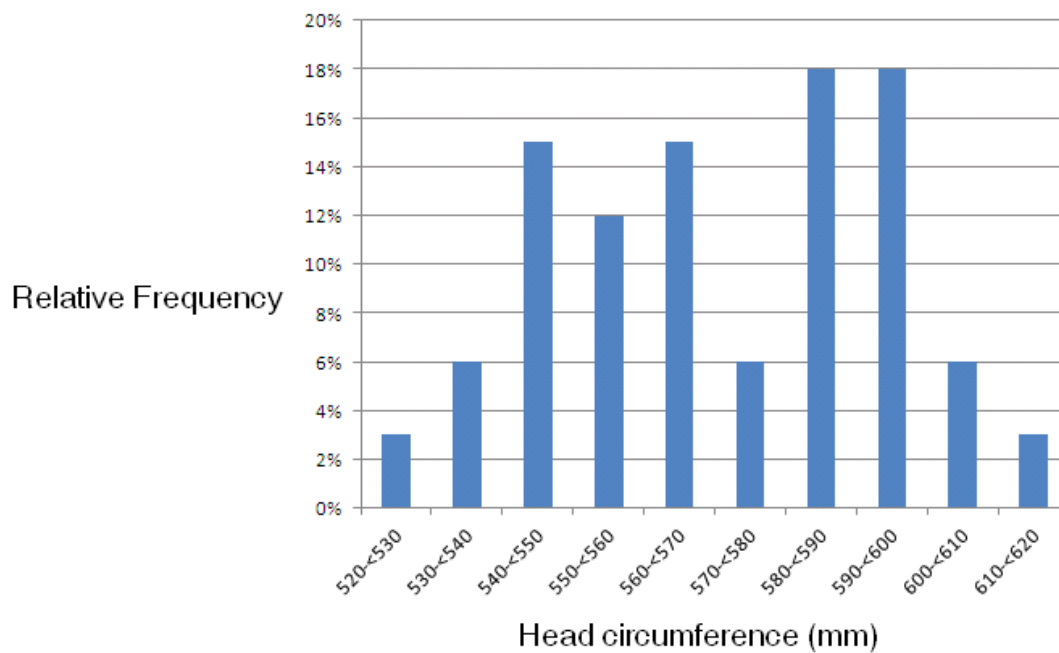
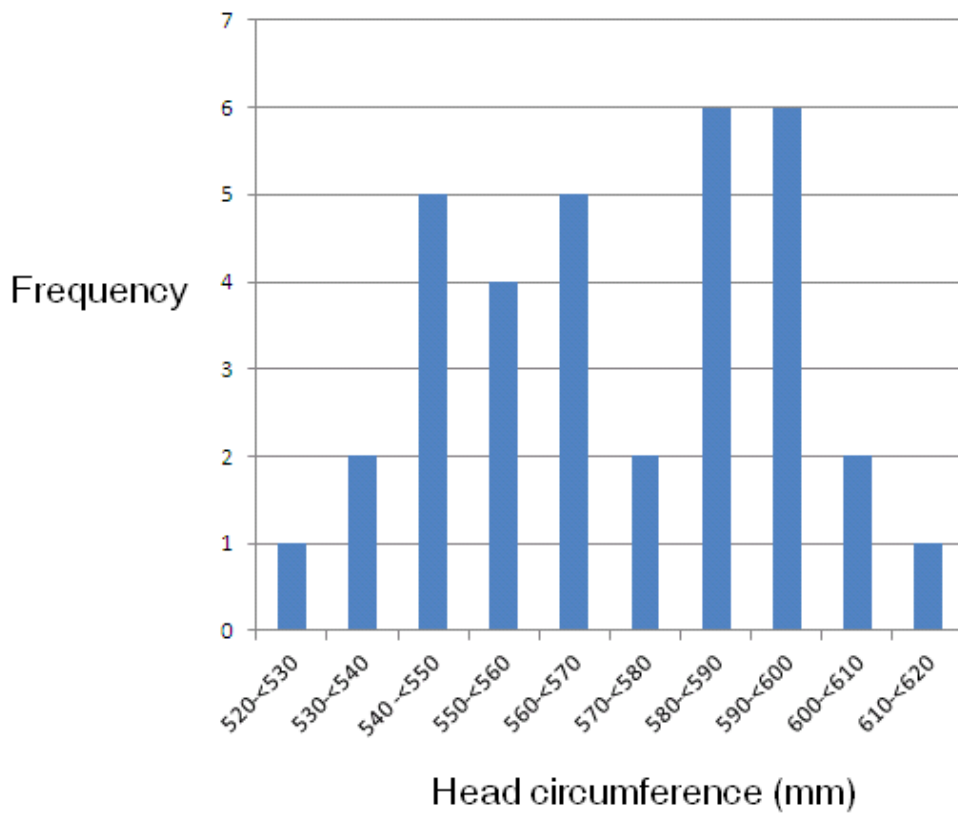
So far you have looked at **line plots** and **stem and leaf plots**. Both are very useful representations for allowing you to see patterns in the variation of your data. A histogram is another useful representation and it is especially useful when dealing with lots of data.

Consider the following:

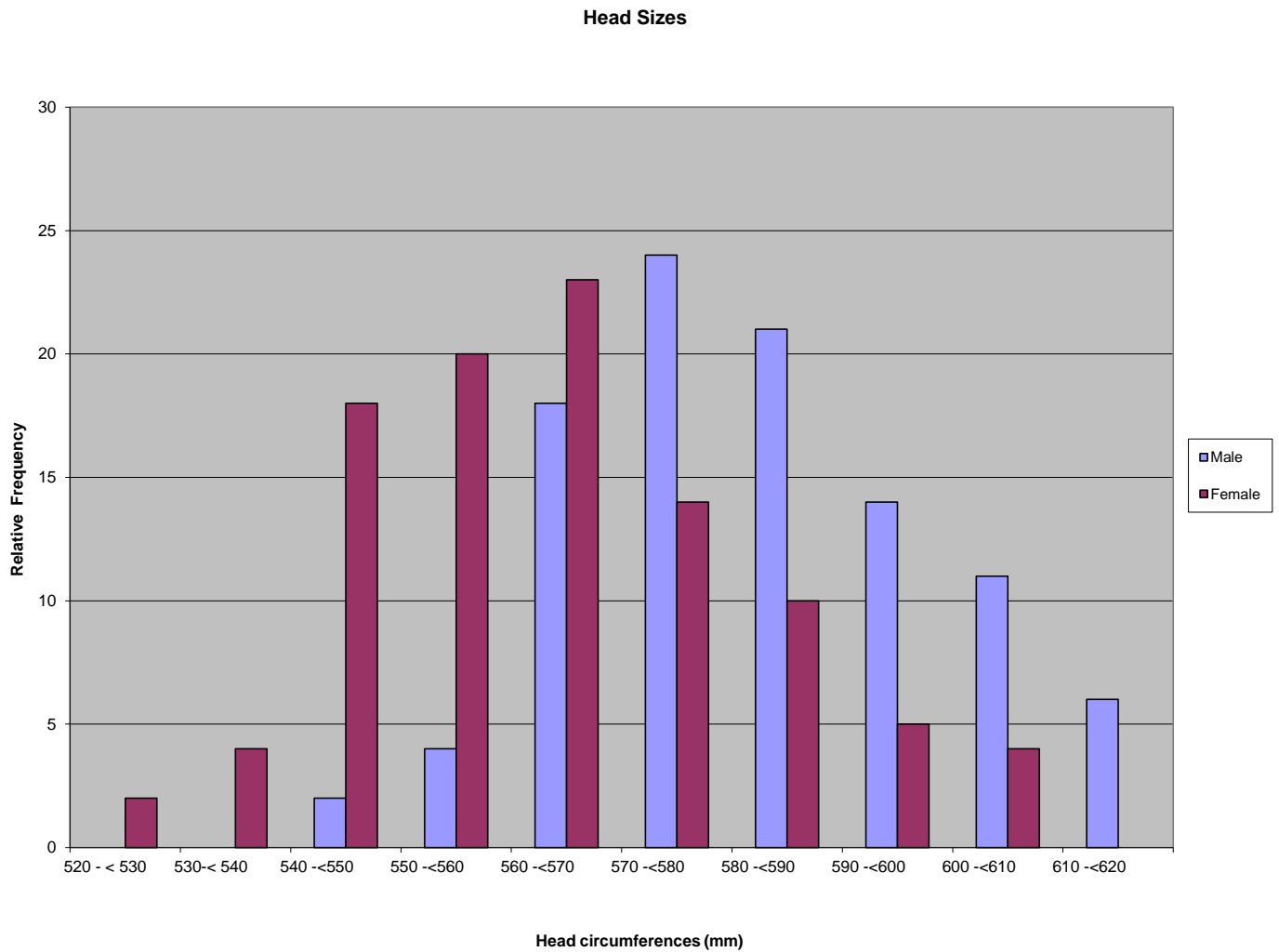
The frequency and relative frequency for each stem was calculated.

			Frequency	Relative Frequency
52	2	520 - < 530	1	1/34 = 3%
53	1 6	530 - < 540	2	2/34 = 6%
54	0 5 5 5 6	540 - < 550	5	5/34 = 15%
55	0 1 2 5	550 - < 560	4	4/34 = 12%
56	0 3 4 5 8	560 - < 570	5	5/34 = 15%
57	3 7	570 - < 580	2	2/34 = 6%
58	0 0 4 5 5 6	580 - < 590	6	6/34 = 18%
59	0 0 0 3 4 5	590 - < 600	6	6/34 = 18%
60	0 0	600 - < 610	2	2/34 = 6%
61	3	610 - < 620	1	1/34 = 3%

Using Excel we can draw a histogram. The diagrams below show two representations. Examine the axes. When would it be more suitable to use relative frequency as opposed to frequency ?

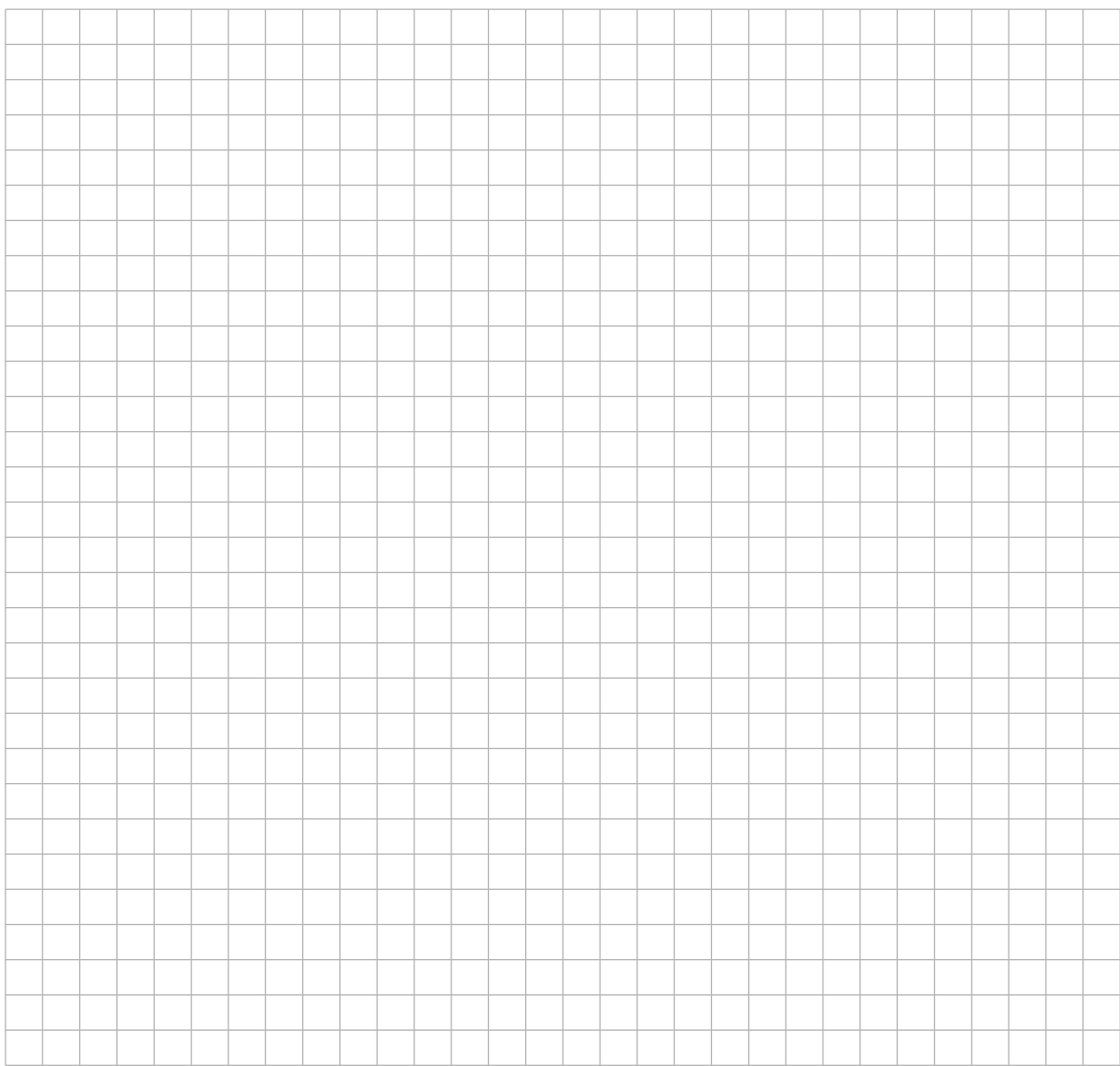


Look at the following histogram showing the distribution of head sizes for a different group of males and females. Compare the distributions. Is there any evidence to suggest that there are differences in the head sizes of men and women?



Why do you think the relative frequency is used for this histogram? Does it matter that the actual numbers of males and females in this sample are not given?

Display this data in a way that will allow you to see patterns in the variation and compare the two **distributions**.



Describe and compare the **shape** of both **sample distributions**.

