

NORMANHURST BOYS HIGH SCHOOL

MATHEMATICS ADVANCED MATHEMATICS EXTENSION 1/2 (YEAR 12 COURSE)



Name:

Initial version by M. Ho, with additional suggestions from H. Lam, June 2020. Last updated November 24, 2021. Various corrections by students and members of the Mathematics Department at Normanhurst Boys High School.

Acknowledgements Pictograms in this document are a derivative of the work originally by Freepik at http://www.flaticon.com, used under CC BY 2.0.

Symbols used

() Beware! Heed warning.	MA12-8 solves problems using appropriate statistical pro-
A Mathematics Advanced content.	cesses
(x) Mathematics Extension 1 content.	Syllabus subtopics
Literacy: note new word/phrase.	MA-F1 Working with Functions
\mathbb{R} the set of real numbers	${\bf MA-S1}~$ Probability and Discrete Probability Distributions
\forall for all	MA-S2 Descriptive Statistics and Bivariate Data Analysis

Syllabus outcomes addressed

Gentle reminder

- For a thorough understanding of the topic, *every* question in this handout is to be completed!
- Additional questions from *CambridgeMATHS Year 12 Advanced* or *Cambridge-MATHS Year 12 Extension 1* will be completed at the discretion of your teacher.
- Remember to copy the question into your exercise book!

Contents

1	Des	cribing Data	4
	1.1	Review of random variables	4
	1.2	Review of types of data	5
	1.3	Review of measures of central tendency	6
	1.4	Review of measures of spread	8
2	Uni	variate Data	10
	2.1	Review of frequency tables	10
		2.1.1 Calculating measures using technology	13
	2.2	Review of frequency histograms and frequency polygons	15
	2.3	(\mathbf{R}) Distributions of data	21
	2.4	Pareto charts	23
	2.5	Review of quartiles and interquartile range	29
	2.6	Review of box plots	31
3	Biva	ariate Data	36
	3.1	Review of two-way tables	36
	3.2	Linear Association of Scatterplots	39
	3.3	Lines of Best Fit	47
4	Priv	vacy, Bias and Ethics	57

Section 1

Describing Data

Learning Goal(s)

E Knowledge

Differentiate between ordinal and nominal data, as well as discrete and continuous data **Calculate the mean, median and standard deviation of data sets**

Vunderstanding

Recognise that discrete data can occur in decimal increments

Solution By the end of this section am I able to:

31.1 Classify data relating to a single random variable

31.3 Calculate measures of central tendency and spread and investigate their suitability in real-world contexts and use to compare large datasets

1.1 (R) Random variables

A random experiment is an experiment with <u>more</u> than <u>one</u> possible outcome.

Definition 2

A random variable X is the <u>outcome</u> of a trial of a random experiment. The various outcomes of the experiment are represented by the values of X.

1.2 (R) Types of data

There are two basic types of data: categorical and numerical.

Definition 3

Categorical data is qualitative and can be grouped in categories. There are two types: 5

1. Ordinal data can be logically ranked in some sort of order.

2. Nominal data has no special order.

Definition 4

 Numerical data is
 quantitative
 and can be counted or measured.

 1. Discrete data has a only take particular values.
 countable only take particular values.
 number of distinct values and can be counted or measured.

2. Continuous data has an <u>infinite</u> number of possible values in a particular range.

Example 1

Determine whether the data sets are nominal, ordinal, discrete or continuous. (a) The weight of fruit taken from each individual tree in an orchard.

- (b) The starting letter or digit of the numberplate for each vehicle in a car park.
- (c) The number of brands of clothing available in a shopping centre.
- (d) The degree of support for the new jumper design for your local sporting team.

								1		REVI	EW OF	MEASU	RES	OF: C	JEN I R	AL I	END	31101
	· · ·		• •			• •	•				• • • • • •				-	• · · ·	•	
		<u> </u>	_		_			_			• •							
1	.3 (F	() (/leas	ures	s of c	entral	ten	dency	y									
		· .	efinit									: : :	: :				:	
	· · · · · · · · · · · · · · · · · · ·										•						•	
	TTI .					verage		C 1			\sum	$x_i f_i$						
÷.	1 ne	mea	$n \frac{x}{\dots}$	or	a	verage	: :	orad	ata se	et:	· · · · · ·	\overline{n}	: •. • :					
	wher	$e x_i$	are t	he .	SCOI	es	, f_i a	re thei	r		frequ	encies			, a	nd r	i is	
					ie samj													
44	une	51	26	OI UI	le sam	JIE.												
÷										•••••								
:) D	efinit	ion 6									: :				:	
	The	med	ian .	Q_2	(seco	nd qua	rtile)	is the		mide	lle	sco	ore	of t	he d	ata	set	
••••	when	arr	angeo	lin		ascendi	no	0	rder									
÷	and the second second																	
	•••••••••••••••••••••••••••••••••••••••	For	an o	ld nu	mber o	of score	s, the	e media	an is t	the	n	niddle		S	core.			
<u>.</u>	•	For	an ey	ven ni	umber	of score	es, the	e medi	an is	the		averag	e		of t	he t	WO	
							, , ,		- 10				••••	••••	0			
-	4		mide	lle	sco	ores												
			·····															
) D	efinit	ion 7														
	The	mod	e is tl	he sco	ore wit	h the		greates	st			freque	ency	7				
														••••	• • • • •			
		A u	nimo	dal da	ata set	has on	е	uniq	ue	. m	ode							
. į	1. A																	
1.0		A	data	set	with	two	or	more	mo	des	is		bi	mod	lal		· · · · ,	
÷.	•					ı two					is	••••	bi	moc	lal	• • •	,	
• • •	•					u two or					is 	••••	bi	moc	lal	••••	,	
											is 		bi	moc	lal	• • •	,	
											is 		bi	moc	lal		,	
											is ···		bi	moc	lal		,	
••••											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											ìs 		bi	moc	lal		,	
											is		bi	moc	lal		,	
											is		bi	moc	lal		,	
											is 		bi	moc	lal		7	
											is 		bi	moc	lal		7	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc			,	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc	lal		,	
											is 		bi	moc			,	
											is 		bi	moc	lal		,	

[2013 General HSC Q14] The July sale prices for properties in a suburb were:

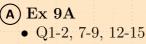
552000, 595000, 607000, 607000, 682000 and 685000.

On 1 August, another property in the same suburb was sold for over one million dollars.

If this property had been sold in July, what effect would it have had on the mean and median sale prices for July?

- (A) Both the mean and the median would have changed.
- (B) Neither the mean nor the median would have changed.
- (C) The mean would have changed and the median would have stayed the same.
- (D) The mean would have stayed the same and the median would have changed.

 $\frac{1}{3}$ Further exercises



x1) Ex 15A

• Q1-2, 7-9, 12-15

NORMANHURST BOYS' HIGH SCHOOL

1.4 (R) Measures of spread

The *interquartile range* is also a measure of spread, which will be covered later.

Definition 8

8

The range of a data set is the difference between the minimum and maximum scores.

Definition 9

of a data set is the average of the The variance σ^2 squared distance of the scores from the mean:

$$\frac{\sum (x_i - \overline{x})^2 f_i}{n} \qquad \text{or} \qquad \frac{\sum x_i^2 f_i}{n} - \overline{x}^2$$

where x_i are the <u>scores</u>, f_i are their <u>frequencies</u>, and n is the total number of scores .

Definition 10

The *standard deviation* σ is a measure of the typical spread of scores from the mean.

The standard deviation of a data set is the square root of the variance:

 $\sqrt{\sigma^2}$

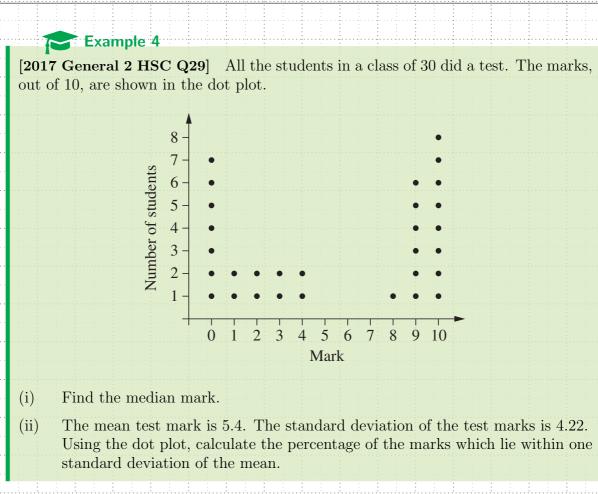
Important note

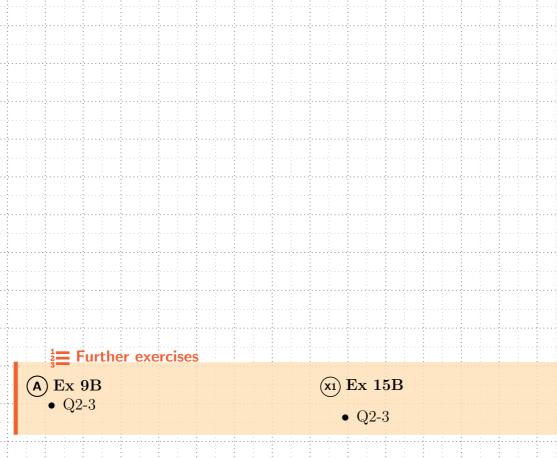
The standard deviation has the same units as the <u>scores</u> in the data set, while the variance has the square units of the scores .

Example 3

[2014 General 2 HSC Q30] The expenditures per primary school student for 15 countries are:

5.9, 7, 7.6, 8.4, 11.2, 11.2, 13.7, 17.1, 18.7, 21.1, 22, 22.5, 23.2, 24.9, 27.6 Calculate the mean, \overline{x} , and the standard deviation, σ , of the data, to two decimal places.





NORMANHURST BOYS' HIGH SCHOOL

Section 2

Univariate Data

Learning Goal(s)

Knowledge

Interpret and construct tabular and graphical displays of data

& Skills Calculate measures of data using a NESA approved calculator to compare data displays

Vunderstanding

Regonise that Pareto charts accentuate the data values with most significant impact

By the end of this section am I able to:

31.2 Organise, interpret and display data into appropriate tabular and/or graphical representations including Pareto charts, cumulative frequency distribution tables or graphs, parallel box-plots and two-way tables

31.4 Summarise and interpret grouped and ungrouped data through appropriate graphs and summary statistics

31.5 Identify outliers and investigate and describe the effect of outliers on summary statistics

31.6 Describe, compare and interpret the distributions of graphical displays and/or numerical datasets and report findings in a systematic and concise manner.

2.1 (R) Frequency tables

Definition 11

A frequency table summarises data and shows frequency values $f_{...}$ for individual or grouped data.

Important note

Grouped numerical data in frequency distribution tables may represent discrete or continuous scores.

Definition 12

Grouping data organising scores into intervals of equal length to provide a clearer overview.

Important note

Grouping data involves ignoring information, which results in its summary statistics to be an approximation of that of the raw data

Definition 13

The class centre (abbreviated to <u>c.c.</u>) is the midpoint of each interval used in the grouping.

Definition 14

The *cumulative frequency* (abbreviated to $\dots \underline{c}, \underline{f}, \dots$) is the number of scores that are less than or equal to a given score, for numerical data.

Important note

Example 5

Number of CDs

0 - 9

10 - 19

20 - 29

30 - 3940 - 49 f

3

5

6 2

1

A frequency distribution table can be extended to a cumulative frequency distribution table by taking the accumulating sums of the frequencies.

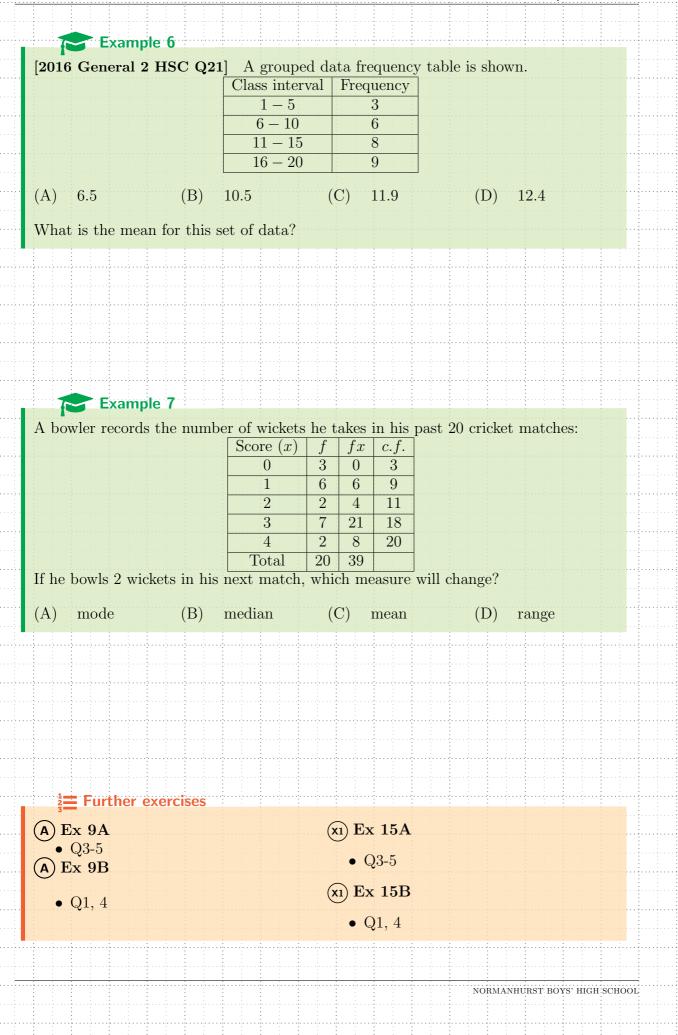
For each of the following frequency, distribution tables, classify the data as discrete

(b)

(a)

Distance travelled (km)	
0 - < 1	3
1 - < 2	5
2 - < 3	6
3 - < 4	2
4 - < 5	1

or continuous:



2	DX/II	7337	OF	FDE	OTT	DIC	v·m	ADT	TR
τı	EVIE	- V.V	Or.	FRE	NO UI	LINC	Ι.Ι.	ADI	പപാ

1 1										1						
7	<u></u>	<u> </u>								· · · · · · · · · · · · · · · · · · ·						
					sures	usin	g teo	chnolo	gy							
			tant i													
All CA	ins SIO <i>f</i>		ions AU F	for PLUS		culato	or fi	unction	IS C	letaile	ed a	are	based	d of	f tl	he
011	510 J	<i>w</i> 02	110 1	LUD		-										
	:= S	tens														
1 - 1 -				into	a frec	uency	v dist	ributio	n tał	ole on	a CA	ASIO	calcu	ilator:		
1.			ey colu			- 	.C ,	SHIF		MODE	, ,	\downarrow ,		for SI		
		1	for [אר (זיז	how (20 row	<u>, ()</u>									
	• • • • • • • • • • • • •	2	tor C)FF (max.	40 ro	ws)									
2.	Ente	er ST	AT mo	ode:	AC	, M	DDE	, 2	for S	STAT ,	1	for 1	L-VAR	1		
3.	Fill	in th	e freq	uenc	y dist	tribut	ion ta	able: us	se the	e arro	w key	rs to 1	naviga	ate th	e cell	s
					· · · · · · · · · · · · · · · · · · ·	nd =		enter o			.		Ŭ			
	:= S	tons							•••••							
	edit th			ev di	stribu	ution t	able.									
•			SHIF1		1,				L fo	or 1-V	AR D	ata				
To r	:eset t	he fr	oquor													
(i)		IIC II	equen	icy u	istrib	ution	table):								
								e: MODE	, [_2	foi	STA	T, 1	fo	r 1-V/	AR	
						, AC			, 2	foi	· STA	T , 🚺	fo	r 1-V/	AR	
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	, 2							
	Ent	er S	FAT m	node	again		; , [, [2 1] ,	e for		T, 1 Edit,		r 1-V <i>I</i> for D		
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							
OR	Ent	er S	FAT m	node	again	, A C	; , [MODE	·							

14	Review of frequency tabi	ES	
: : :	E Steps		
÷.	Calculating 'Variance': AC , SHIFT 1 , 4 for Var, then:		•
	• mean: $\begin{vmatrix} 2 \\ 0 \end{vmatrix}$ for \overline{x} , $\begin{vmatrix} = \\ 0 \end{vmatrix}$		
	• standard deviation: 3 for $\sigma \mathbf{x}$, =		
	• size: 1 for n , =		
	📰 Steps		
	• sum of scores: 2 for $\Sigma \mathbf{x}$, =		
	• sum of squared scores: press 1 for Σx^2 , =		
	• sum of squarea scores. press 1 for Δx , -		
: 			
	📰 Steps		
	• <i>minimum</i> : 1 for minX , =		
	• maximum: 2 for maxX, =		
			• • • • • •
	• lower quartile: 3 for Q_1 , =		
			• • • • • •
	• median: 4 for med , =		
÷.	• upper quartile: 5 for Q_3 , =		
		• • • • • • • • •	•
	Important note		
	The median, lower quartile and upper quartile functions are only available on		
	CASIO fx-82 PLUS II model calculators and above.		
N 8 8 8 1 1 1 1 1 1			
			•
	NORMANHURST BOYS' HIGH SCHO	OL	

2.2 (R) Frequency histograms and frequency polygons

A frequency histogram is visually similar to a <u>column</u> graph, but with no <u>gaps</u> in between that columns.

The subintervals on the horizontal axis of frequency histograms are often called *bins*.

Important note

When constructing a frequency histogram:

- The first column is usually placed one half-column width from the vertical axis.
- The columns join up with no gaps.
- Each column is centred on the value for individual data.
- Each column is centred on the class centre for grouped data.
- A *scale break* or zigzag on the x or y-axis indicates that the data displayed does not include all the values that exist on the number line used.

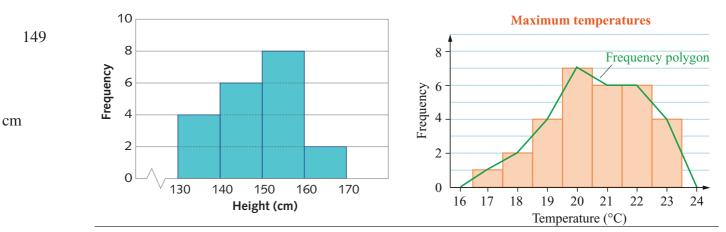
Definition 16

A *frequency polygon* is a <u>line</u> <u>graph</u> often constructed on the same graph as a frequency histogram.

Important note

When constructing a frequency polygon:

- The plotted points are at the centre of the top of each column of the histogram.
- The plotted points are joined with intervals.
- The polygon starts on the left, on the horizontal axis at the previous value or class centre.
- The polygon ends on the right, on the horizontal axis at the next value or class centre.



Important note

Coarser grouping in frequency histograms is more practical, as too many columns can make the data display difficult to interpret.

Example 8

At the start of Year 7, Cedar Heights High School gave 40 students a spelling test marked out of 10. The raw results were organised into a frequency table.

Mark (x)	1	2	3	4	5	6	7	8	9	10	
Frequency (f)	2	4	2	1	6	8	7	6	2	2	

- (a) Draw a histogram and frequency polygon for the original data.
- (b) Construct a frequency distribution table by grouping the data into subintervals of 1-2, 3-4, ..., including a row for the class centres.
- (c) Draw a histogram and frequency distribution polygon for the grouped data.
- (d) Compare and comment on what the two data displays have shown.

NORMANHURST BOYS' HIGH SCHOOL

Definition 17

A *cumulative frequency histogram* is formed by stacking the columns of a frequency histogram.

Important note

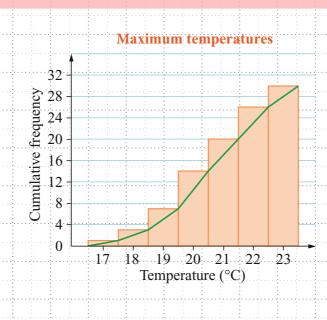
Definition 18

A cumulative frequency polygon or ogive shows cumulative frequency and so is drawn slightly differently.

Important note

When constructing a cumulative frequency polygon:

- The polygon starts at zero, at the bottom left-hand corner of the first column (no scores have yet been accumulated).
- The polygon passess through the top right-hand corner of each column (scores less than or equal to the upper bound of the class interval are plotted).
- The polygon finishes at the top right-hand corner of the last column (its height equals the total size of the sample),



NORMANHURST BOYS' HIGH SCHOOL

Returning to the spelling test marks from Example 8:

Mark (x)	1	2	3	4	5	6	7	8	9	10
Frequency (f)	2	4	2	1	6	8	7	6	2	2

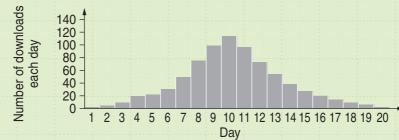
- (a) Insert a cumulative frequency row and construct a cumulative frequency histogram and ogive for the data.
- (b) Group the data by pairing the marks and construct a grouped frequency distribution table, including a cumulative frequency row.

(c) Construct a cumulative frequency histogram and ogive for the grouped data.

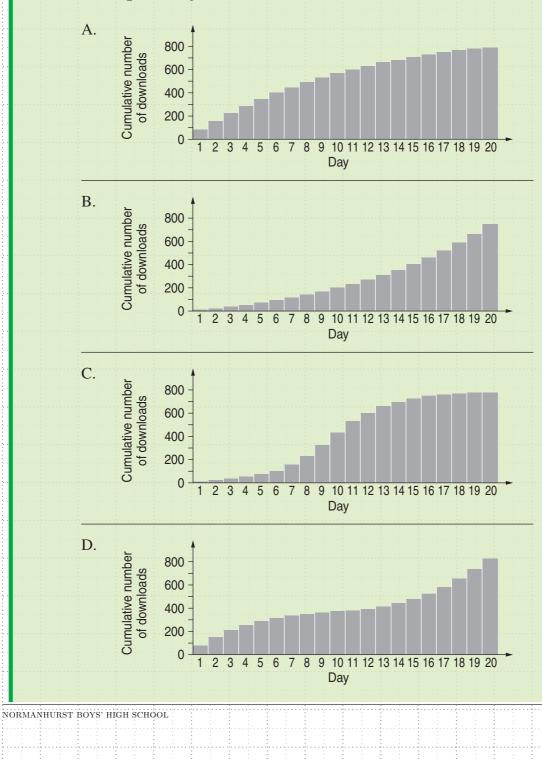
(d) Find the median of the original data and the grouped data, and compare them.

NORMANHURST BOYS' HIGH SCHOOL

[2021 Adv HSC Q4] The number of downloads of a song on each of twenty consecutive days is shown in the following graph.



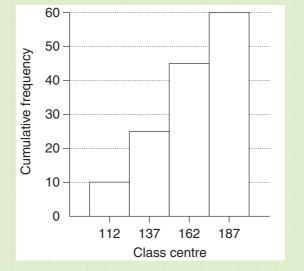
Which of the following graphs best shows the cumulative number of downloads up to and including each day?



[2010 General HSC Q26] A new shopping centre has opened near a primary school. A survey is conducted to determine the number of motor vehicles that pass the school each afternoon between 2:30 pm and 4:00 pm.

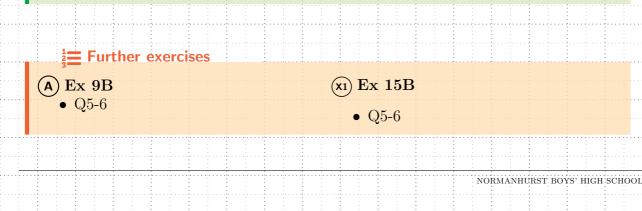
The results for 60 days have been recorded in the table and are displayed in the cumulative frequency histogram.

 Score	Class centre	Frequency	Cumulative frequency
100 - 124	112	10	10
 125 - 149	137	×	25
150 - 174	162	20	45
 175 - 199	187	15	60



- (i) Find the value of \times in the table.
- (ii) Carefully copy the cumulative frequency histogram and draw a cumulative frequency polygon (ogive) for this data.
- (iii) Use your graph to determine the median. Show, by drawing lines on your graph, how you arrived at your answer.
- (iv) Prior to the opening of the new shopping centre, the median number of motor vehicles passing the school between 2 : 30 pm and 4 : 00 pm was 57 vehicles per day.

What problem could arise from the change in the median number of motor vehicles passing the school before and after the opening of the new shopping centre? Briefly recommend a solution to this problem.



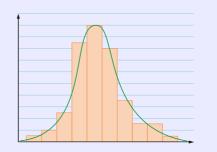
2.3 (R) Distributions of data

The shape of a statistical display indicates the distribution of the data.

Definition 19

Clustering occurs when the scores are close together or 'bunched up'.

Definition 20

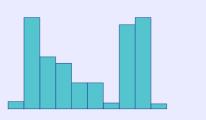


Important note

Distributions with scores that are approximately evenly spread about either side of the centre of distribution, are still classified as symmetrical distributions.

Definition 21

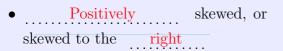
The distribution is <u>bimodal</u> as it has <u>two</u> modes.



ICE-EM Mathematics 9 3ed ISBN 978-1-108-40432-7 © The University of Melbourne / AMSI 2017 Cambridge University Press Photocopying is restricted under law and this material must not be transferred to another party.

Definition 22





• Most of the scores are relatively <u>low</u>...



- <u>Negatively</u> skewed, or skewed to the <u>left</u>
- Most of the scores are relatively high

2.4 Pareto charts

Definition 23

• The *Pareto chart* consists of a column graph with columns arranged in <u>descending</u> order and a <u>cumulative</u> percentage line graph, drawn together on the same chart.

Fill in the spaces

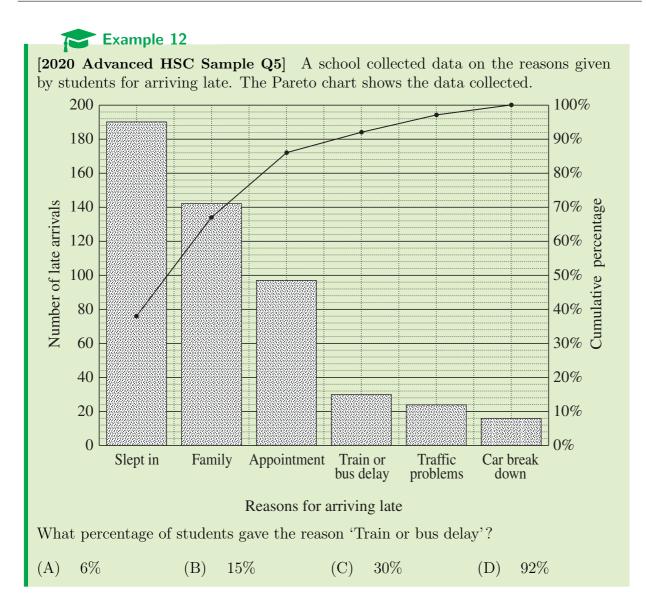
- The chart usually has two vertical axes: frequencies on the left and percentage frequencies on the right.
- The purpose of the chart is to identify the most significant significant

The line graph indicates the added contribution of each problem.

- If the line graph rises steeply and then levels out, the first two or three problems have the <u>most</u> impact.....
- If the line graph rises at a steady rate, all problems have roughly equal impact and so the columns will be at similar heights

Important note

- The column graph is plotted using the frequency column of the frequency table, arranged in descending order.
- The cumulative line graph is plotted using a cumulative percentage column.



An online seller of clothing summarised the complaints received in a month in the following table.

Type of complaint	No. of complaints	Percentage	Cumulative Percentage
Problems completing the	50		
order online			
Difficulty accessing website	30		
Cancelled order	9		
Wrong article sent	5		
Overcharging for delivery	4		
Late delivery	2		

(a) Fill in the table's *percentage* and *cumulative percentage* columns.

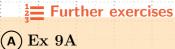
(b) Construct a Pareto chart for this information.

(c) Which problems account for 80% of the complaints?

A bank wants to identify the most significant problems it experiences in the delay of processing credit card applications. A random survey was conducted and the results are shown in the table.

	Cause of delay	Frequency (f)
	Incorrect address	36
	Can't read	17
	No signature	73
•	Wrong from	11
	Other	3

- (a) Arrange the problems by frequency in descending order and insert a cumulative percentage column.
- (b) Construct a Pareto chart of the survey data.
- (c) What area(s) should the manager concentrate on improving?
- (d) What percentage of her problem does this solve?



• Q9-11





NORMANHURST BOYS' HIGH SCHOOL

Exercises Source: (?, ?, Ex 7G Q6-9)

P	6	A random survey of customer complaints yielded the data shown	Type of complaint	Freque
ROB		on the right.	Packaging	11
LEM		a Draw a Pareto chart to illustrate this information.b What area(s) should this company concentrate on	Delivery	7
2 O L		improving?	Invoice wrong	36
VING		c What percentage of the problems would this solve?	Product quality	22
PROBLEM SOLVING, REASONING AND JUSTIFICATION			Other	4
I N O S A	7	The owner of a shoe store takes a random sample of customer		
NG	-	complaints. The results are shown in the table on the right.	Type of complaint	Freque
A N D		a Draw a Pareto chart to illustrate this information.	Difficult parking	77
JU		b Before the survey, the manager thought that it was the limited	Salesperson rude	9
STIF		range of styles being offered that was the main reason for the	Poor lighting	5
I C A I		decline in her business and she blamed the supplier. What	Layout confusing	8
1 O N		percentage of the problem was caused by limited styles?	Limited sizes	37
		c If you were the shoe store owner, what area(s) would you concentrate on improving?	Limited styles	11
		d What percentage of the problems would your improvements	Other	3
	8	 Pareto charts can be drawn using a spreadsheet. After entering the d Insert > Chart > Histogram > Pareto. a Use your computer spreadsheet to draw Pareto charts for question b What is different about the charts drawn by Excel compared to y 	ons 5–7 .	ole, go to:
CH	9	A restaurant manager takes a random survey of customer	Type of complaint	Freque
CHALLENGE		complaints in order to increase the patronage of his restaurant.	Rude staff	6
ENG		The results are shown in the table on the right.a Draw a Pareto chart to illustrate this information.	No atmosphere	8
m		b If you were the manager, what area(s) would you concentrate	Small portions	17
		on improving?	Too noisy	19
		c What percentage of the problems would your improvements	Too expensive	93
		from part b solve?	Limited menu	5
			Dirty washrooms	8 35
			Long delays in serving	55

6 a	Type of complaint	Frequency	Relative frequency (%)	Cumulative relative frequency (%)
	Invoice wrong	36	45	45
	Product quality	22	27.5	72.5
	Packaging	11	13.8	86.3
	Delivery	7	8.7	95
	Other	4	5	100
	Total	80		

S
2
ш
≥
S
Z

◄

Type of complaint	Frequency
Packaging	11
Delivery	7
Invoice wrong	36
Product quality	22
Other	4

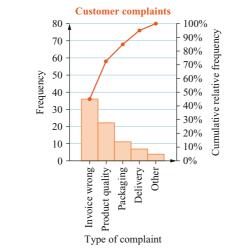
Type of complaint	Frequency
Difficult parking	77
Salesperson rude	9
Poor lighting	5
Layout confusing	8
Limited sizes	37
Limited styles	11
Other	3

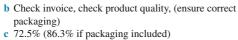
9

Cramped seating

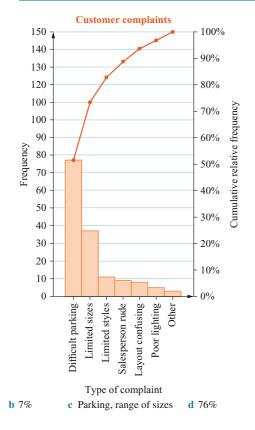
S ۲ ш NSN ∢

7



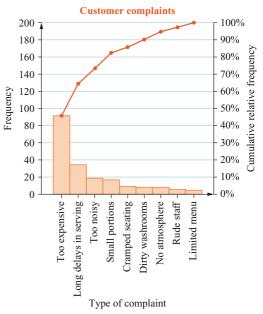


a	Type of complaint	Frequency	Relative frequency (%)	Cumulative relative frequency (%)		
	Difficult parking	77	51	51		
	Limited sizes	37	25	76		
	Limited styles	11	7	83		
	Salesperson rude	9	6	89		
	Layout confusing	8	5	94		
	Poor lighting	5	3	97		
	Other	3	2	100		
	Total	150				



a	Type of complaint	Frequency	Relative frequency (%)	Cumulative relative frequency (%)				
	Too expensive	93	46.5	46.5				
	Long delays in serving	35	17.5	64				
	Too noisy	19	9.5	73.5				
	Small portions	17	8.5	82				
	Cramped seating	9	4.5	86.5				
	Dirty washrooms	8	4	90.5				
	No atmosphere	8	4	94.5				
	Rude staff	6	3	97.5				
	Limited menu	5	2.5	100				
	Total	200						

9



b Cost of meals, delays in serving, noise, (portion size)

c 74% (82% if portion size included)

2.5 (R) Quartiles and interquartile range

A data set can be divided into <u>four</u> parts by <u>three</u> quartiles; the lower quartile, median and upper quartile.

The *lower quartile* Q_1 is the <u>median</u> of the lower <u>half</u> of the data set.

The upper quartile Q_3 is the median of the upper half of the data set.

Important note

If the data set has an odd number of scores, the median is excluded when finding the lower and upper quartiles.

Definition 25

The *interquartile range* IQR measures the spread of the middle 50% of the data set:

$$IQR = Q_3 - Q_1$$

Important note

The interquartile range is often a better measure of the spread of the data than the range, particularly when the data set contains outliers.

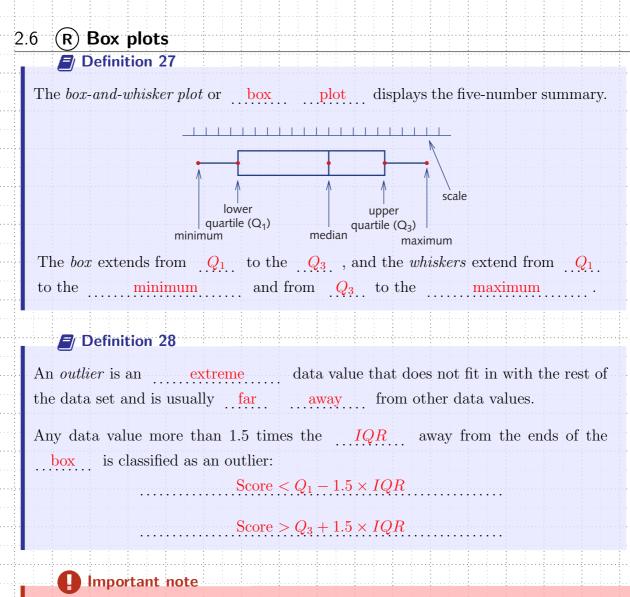
Definition 26

The *five-number summary* is a useful summary of a data set:

- **1.** Minimum
- 2. Lower quartile
- 3. Median
- 4. Upper quartile
- 5. Maximum

D .								
Review	OF	QUARTILES	AND	INTERC)UA	BTILI	$\mathbf{E} \cdot \mathbf{R} \mathbf{A}$	NGE

2016 Ge		HSC					e wr	ote d	lown	the	num	ber of	fgoal	s sco	red	
ı 9 differ	rent gai	mes du	ring t		eason. $3, 3,$		5, 8,	9. □								
he last i	number	has be	een or							a is 1	.0.					
Vhat is t	the five	-numbe	er sun	nmar	y for	this d	ata	set?								
						•						• • • • • • • • • • • • • • • • • • •				
											••••					
	* *											· · · · · · · · · · · · · · · · · · ·				
						•						• • • • • • • • • • • • • • • • • • •				
												•				
		·····										· · · · · · · · · · · · · · · · · · ·				
						• • • • • • • • • • • • • • • • • • •						• • • • • • • • • • • • • • • • • • •				
		······				••••••••••••••••••••••••••••••••••••••						• • • • • • • • • • • • • • • • • • •	(
						· · · · · · · · · · · · · · · · · · ·						• • • • • • • • • • • •				
												• • • • • • • • • • • • • • • • • • •				
				•••••												•••••
			· · · · ·			•			-				1			
	Exam	ole 16														
a se 📕 de la serie	Examp meral 2	2 HSC	Sc	ore		1	2	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3			6	data	set is	shov	wn.	
2018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	sho	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
2018 Ge Zhat is t	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	
018 Ge	eneral 2	2 HSC	Sc ulativ	core re fre	quenc	1 y 5	29	3	4	5	6	data	set is	shov	wn.	



When constructing a box plot for a data set with outliers, the whiskers end at the highest and lowest data values that lie within $1.5 \times IQR$ from the ends of the box.

The outliers are marked with dots.

Example 17

[2017 General 2 HSC Q30] A set of data has a lower quartile (Q_1) of 10 and an upper quartile (Q_3) of 16.

What is the maximum possible range for this set of data if there are no outliers?

© The University of Melbourne / AMSI 2011

NORMANHURST BOYS' HIGH SCHOOL

32

The waiting times in seconds at a ticket counter were as follows:

 $\begin{matrix} 0, 0, 3, 5, 5, 5, 9, 10, 12, 13, 16, 17, 18, 18, 21, 22, 23, 23, 24, 24, 24, 24, 24, \\ 25, 25, 25, 26, 26, 27, 28, 28, 28, 29, 29, 29, 30, 31, 31, 31, 32, 33, 33, 33, \\ 34, 34, 34, 34, 35, 35, 35, 36, 36, 37, 38, 38, 38, 39, 39, 39, 40, 41, 41, 52 \end{matrix}$

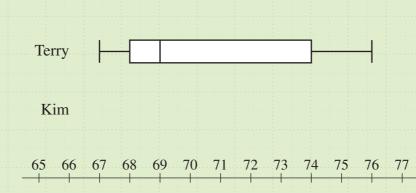
- (a) Find Q_1 , the median, Q_3 and the IQR.
- (b) Find any outliers.
- (c) Draw a boxplot, showing outliers.
- (d) Comment on the shape of the boxplot.

Answer: (a) $Q_1 = 22.5$, median = 29, $Q_3 = 34.5$, IQR = 12 (b) 0, 0, 3 (d) Negatively skewed

[2014 General 2 HSC Q29] Terry and Kim each sat twenty class tests. Terrys results on the tests are displayed in the box-and-whisker plot shown in part (i).

(i) Kims 5-number summary for the tests is 67, 69, 71, 73, 75.

Draw a box-and-whisker plot to display Kims results below that of Terrys results.

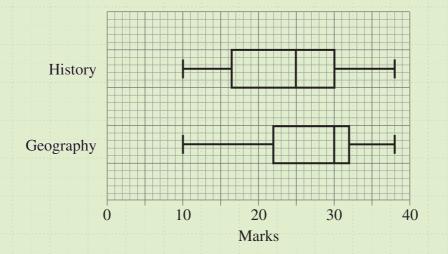


- (ii) What percentage of Terry's results were below 69?
- (iii) Terry claims that his results were better than Kims. Is he correct? Justify your answer by referring to the summary statistics and the skewness of the distributions.



34

[2016 General 2 HSC Q22] The box-and-whisker plots show the results of a History test and a Geography test.



In History, 112 students completed the test. The number of students who scored above 30 marks was the same for the History test and the Geography test.

How many students completed the Geography test?

	. ((\mathbf{A}))	8				((B)	50			(C)	Ę	56			(Γ))	1	12					
••••	····}				 		:			 	 	 					 	 					••••		•••••	•••••	
							: : :																				
••••		•••••			 		: : :			 	 	 					 	 • • • • •									
	•••••																										
••••	•••••	•••••			 •			• • • • •		 	 	 			• • • • •		 	 	•••••			• • • • •			•••••	••••	
••••	•••••	•••••					 			 		 					 	 • • • • •	•••••		•••••		••••		•••••	•••••	
					 				· · · · · · ·	 	 	 					 	 									
;					 · · · · · ·					 	 	 					 	 									
		•••••			 	••••		• • • • • •		 	 	 					 	 									
					•														•								
		•••••			 	• • • • •		• • • • • •		 	 	 					 	 	•••••								
	•••••				•																						
					 			•••••		 	 	 					 	 									
					•																						
	:				 •							 •					 	 	•••••						•••••	•••••	
					•																						
					 					 	 	 					 	 				DOT					
					•														INOI	πМА	NHU	кSТ	ROJ	S H	IGH	SCH	OOL
							:			 	 							 	•							•••••	
									. :																		

The towns of Karuah and Buladelah both have a speed limit of 60 km/h. The speeds of the first 100 cars travelling through these towns, from 9:00 am to 10:00 am on Saturday, were measured and the results are shown in the table.

	Karuah	Buladelah
Mean	59	60
Median	58	58
Lower quartile	52	50
Interquartile range	18	9
Highest speed	85	105
Lowest speed	35	40

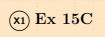
(a) Construct parallel box plots to display this data.

(b) Explain why the police might target Karuah for speeding.

(c) Explain why the police might target Buladelah for speeding.







• Q1-8

NORMANHURST BOYS' HIGH SCHOOL

Section 3

Bivariate Data

Learning Goal(s)

Knowledge

Identify and describe linear association in bivariate scatterplots

🗘 Skills

Calculate Pearson's correlation coefficient and the equation of the least squares regression line using a NESA approved calculator

V Understanding

Recognise the possible but not guaranteed link between correlation and causation

By the end of this section am I able to:

- 31.7 Construct a bivariate scatterplot to identify patterns in the data that suggest the presence of an association
- 31.8 Use bivariate scatterplots (constructing them where needed), to describe the patterns, features and associations of bivariate datasets, justifying any conclusions
- 31.9 Calculate and interpret Pearsons correlation coefficient (r) using technology to quantify the strength of a linear association of a sample
- 31.10 Model a linear relationship by fitting an appropriate line of best fit to a scatterplot and using it to describe and quantify associations
- 31.11 Use the appropriate line of best fit, both found by eye and by applying the equation of the fitted line, to make predictions by either interpolation or extrapolation
- 31.12 Solve problems that involve identifying, analysing and describing associations between two numeric variables
- 31.13 Construct, interpret and analyse scatterplots for bivariate numerical data in practical contexts

Definition 29

A *bivariate data set* consists of two different <u>variables</u> for each data point.

3.1 (R) Two-way tables

A two-way table or contingency table consists of two or more related frequency tables combined together.

Information can be read from each <u>cell</u>, <u>row</u> or <u>column</u>, with each <u>row</u> and <u>column</u> representing a separate frequency table.

[2016 General 2 HSC Q23] A group of 485 people was surveyed. The people were asked whether or not they smoke. The results are recorded in the table.

	Smokers	Non-smokers	Total
Male	88	176	264
Female	68	153	221
	156	329	485

A person is selected at random from the group.

What is the approximate probability that the person selected is a smoker OR is male?

	(A)	-33	3%			((B)		189	%			C)	(58%	6			(Γ))	8	7%)		
																			•							
							 					 	 							 					 	••••
							 												•							
				: : :			 					 : 	 					: :	: :	 		: : :			 	
																			• • •			•				
••••							 					 	 					••••		 					 •••••	• • •
	-	7		? F	yai	mn	23		-										-							



[2011 General HSC Q25] At another school, students who use mobile phones were surveyed. The set of data is shown in the table.

	Pre-paid	Plan	TOTAL
Female students	172	147	319
Male students	158	103	261
TOTAL	330	250	

(i) How many students were surveyed at this school?

(ii) Of the female students surveyed, one is chosen at random. What is the probability that she uses pre-paid?

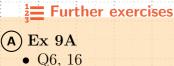
Ten new male students are surveyed and all ten are on a plan. The set of data (iii) is updated to include this information. What percentage of the male students surveyed are now on a plan? Give your answer to the nearest per cent.

[2017 General 2 HSC Q29] A group of Year 12 students was surveyed. The students were asked whether they live in the city or the country. They were also asked if they have ever waterskied.

The results are recorded in the table.

	Have waterskied	Have never waterskied
Live in the city	150	2500
Live in the country	70	800

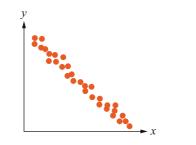
- (i) A person is selected at random from the group surveyed. Calculate the probability that the person lives in the city and has never waterskied.
- (ii) A newspaper article claimed that Year 12 students who live in the country are more likely to have waterskied than those who live in the city. Is this true, based on the survey results? Justify your answer with relevant calculations.

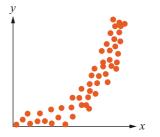






The form of bivariate data sets can be <u>linear</u> or <u>non-linear</u> relationship, depending on the shape of its clustering:





Important note

Linear relationships between two variables can be described in terms of direction and strength of association.

Definition 33

The *direction* of a linear relationship can be negative or positive:

- In a *negative relationship*, one variable increases as the other variable decreases .

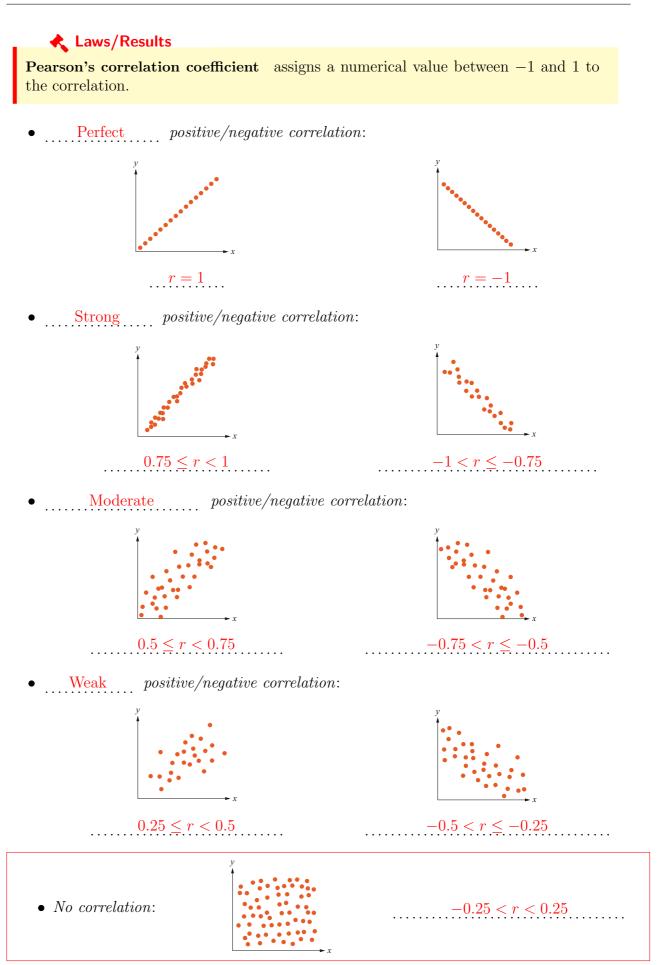
Definition 34

 Pearson's correlation coefficient,
 r.
 is a statistical measure of the strength of

 association
 (or
 correlation
) between two variables.

Important note

The formula for Pearson's correlation coefficient is not developed or used by hand in Mathematics Advanced/Extension 1 or Extension 2.



	🚍 Steps
	calculate Pearson's correlation coefficient on a CASIO calculator:
	Bivariate frequency table: MODE 2 for STAT, 2 for A+BX
2.	Use the arrow and numbers keys to enter the data points into the table.
3.	AC , SHIFT 1 , 5 for Reg , 3 for r , =

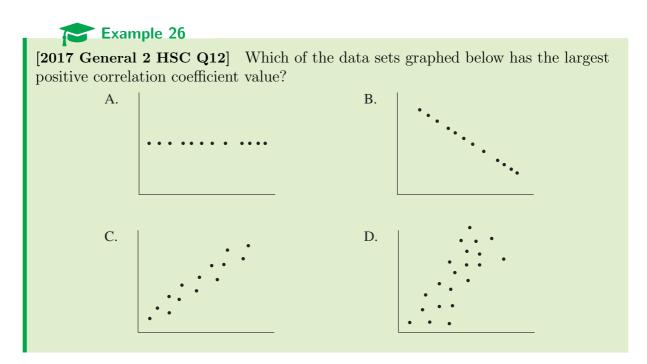
Important note

Correlation does not necessarily imply causation.

Example 25

[2012 General HSC Q11] Which of the following relationships would most likely show a negative correlation?

- (A) The population of a town and the number of hospitals in that town.
- (B) The hours spent training for a race and the time taken to complete the race.
- (C) The price per litre of petrol and the number of people riding bicycles to work.
- (D) The number of pets per household and the number of computers per household.



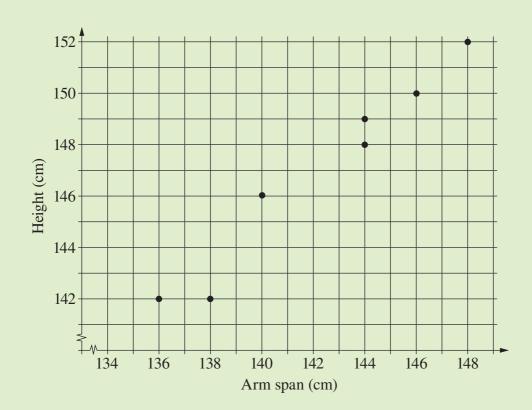


[2010 General HSC Q6] A survey of Year 7 students found a number of relationships with a high degree of correlation.

Which of the following relationships also demonstrates causality?

- (A) Students height and the length of their arm span
- (B) The size of students left feet and the size of their right feet
- (C) Students test scores in Mathematics and their test scores in Music
- (D) The number of hours students spent studying for a test and their results in that test

[2019 Standard 2 HSC Q23] A set of bivariate data is collected by measuring the height and arm span of seven children. The graph shows a scatterplot of these measurements.



- (a) Calculate Pearsons correlation coefficient for the data, correct to two decimal places.
- (b) Identify the direction and the strength of the linear association between height and arm span.
- (c) The equation of the least-squares regression line is shown.

 $\text{Height} = 0.866 \times (\text{arm span}) + 23.7$

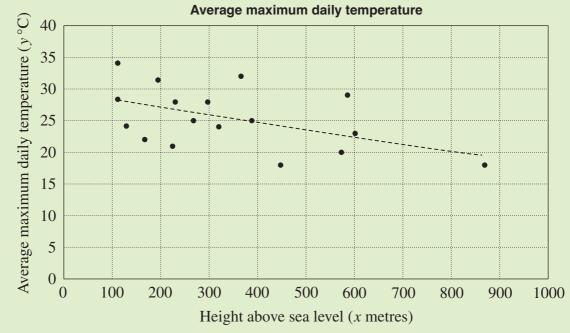
A child has an arm span of 143 cm.

Calculate the predicted height for this child using the equation of the least-squares regression line.

Answer: (a) 0.98 (b) Positive, strong (c) 147.538 cm

[2021 Adv HSC Q17] For a sample of 17 inland towns in Australia, the height above sea level, x (metres), and the average maximum daily temperature, y (°C), were recorded.

The graph shows the data as well as a regression line.



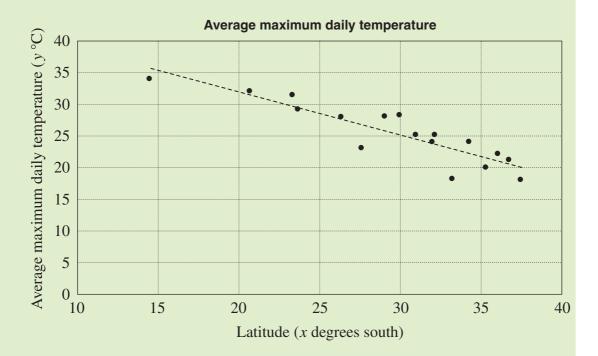
The equation of the regression line is y = 29.2 - 0.011x.

The correlation coefficient is r = -0.494.

- (a) (i) By using the equation of the regression line, predict the average maximum daily temperature, in degrees Celsius, for a town that is 540 m above sea level. Give your answer correct to one decimal place.
 - (ii) The gradient of the regression line is -0.011. Interpret the value of this gradient in the given context.

[2021 Adv HSC Q17] Example 29 on the preceding page continued...

(b) The graph below shows the relationship between the latitude, x (degrees south), and the average maximum daily temperature, y (°C), for the same 17 towns, as well as a regression line.

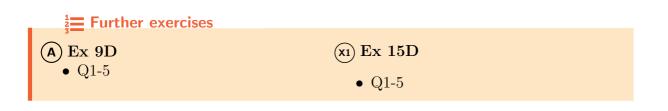


The equation of the regression line is y = 45.6 - 0.683x.

The correlation coefficient is r = -0.897.

Another inland town in Australia is 540 m above sea level. Its latitude is 28 degrees south.

Which measurement, height above sea level or latitude, would be better to use to predict this town's average maximum daily temperature? Give a reason for your answer.



3.3 Lines of Best Fit

The *line of best fit* or <u>trendline</u> is a <u>straight</u> line that provides a representation of all the data points in a <u>scatterplot</u> that has a <u>linear</u> correlation.

Important note

Draw a line of best fit by eye involves fitting the line such that:

- the distance of points from the line is minimised
- an approximatelyequal number of data points lie above and below the line The line does not need to pass through any data points.

Definition 36

The *least-squares regression line* is a mathematically determined straight line that <u>best</u> fits the data set by <u>minimising</u> the squares of the vertical distances from the points to the line:

y = mx + b

where m =________ and b =__________ *y*-intercept

A Laws/Results

The gradient of the least-squares regression line:

$$\dots m = r \times \frac{\sigma_y}{\sigma_x}$$

where $r = $ Pearson's	correlation	coefficient
$\sigma_x = ext{standard}$	deviation	\dots of the <i>x</i> -variable,
$\sigma_y = ext{standard}$	deviation	\dots of the <i>y</i> -variable.

A Laws/Results

The *y*-intercept of the least-squares regression line:

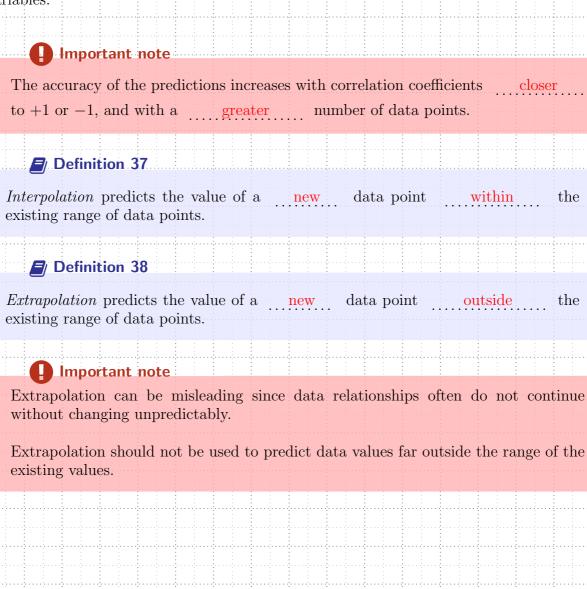
 $b = \overline{y} - m\overline{x}$

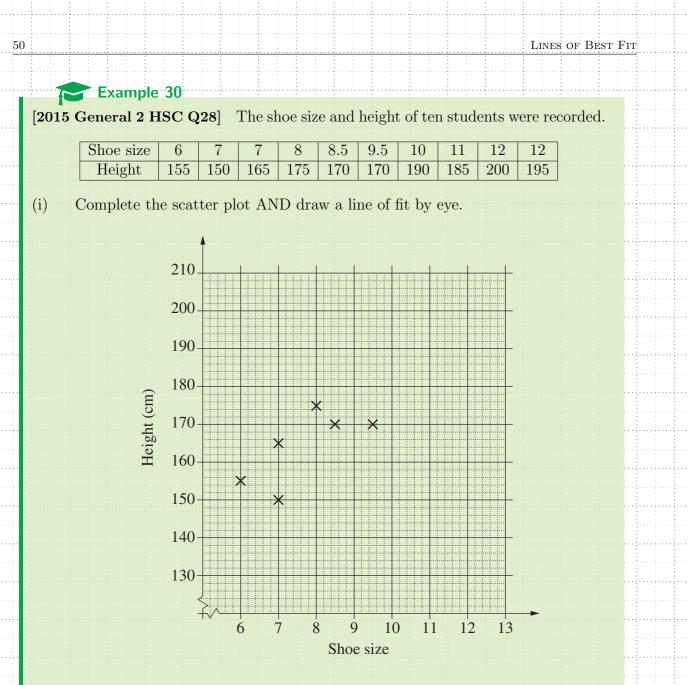
where m = gradient , $\overline{x} =$ mean of the *x*-variable, and $\overline{y} =$ mean of the *y*-variable

48	· · · · · · · · · · · · · · · · · · ·	 					*				· · · · · · · · · · · · · · · · · · ·				· · · · · · · · · · · · · · · · · · ·	· · · · · · · · · · · · · · · · · · ·		*		Lı	VES	of Be	st I	7IT				
· · · · · · ·	•			aws	/Re	esult	5											*										
	Th							t-sq	uare	es re	egres	ssion	lin	e ca	n al	so be	exp	ress	ed a	ls:								
											y =	= A +	- <i>B</i> :	r														
	wh	ere	e A	=	• • • •	<u> </u>	-int	erce								gradi	ient	•••	•									•
																									· · · · · · · · · · · · · · · · · · ·			•
	То	_	nd t			tion	of t	he	least	t-sqi	uare	es reg	gres	sion	ı line	e on a	ı CA	SIC) cal	cula	ato	r:						
- E	1.					frequ							,															
			•	M	ODE		2	for	STA	Γm	ode																	
			٠	2	f	or A	+BX																		•••		• • • •	
1 3	2.						and	nur	nber	rs ke	eys	to er	nter	the	e dat	a poi	nts i	nto	the	tab	ole.							
	3.		y-ir								1 -																	
			•	A			IFT		1	5	for	r Reg	5															
				1		or A	, the	en [• • • • • • • • • • • • • • • • • • • •			•
	4.		gra	-																								• •
			•	A	<u></u>	· <u>· · · · · · · ·</u>	IFT		1	5	tor	r Reg	5												•••		••••	
			•	2		or B.	, the	en [=									:							•••••		••••	
· · · · ·					••••••												• • • • • • •					· · · · · · · · · · · · · · · · · · ·			•••			•
					•••																	· · · · · · · · · · · · · · · · · · ·						
:																						· · · · · · · · · · · · · · · · · · ·						
				· · · · · · · · · · · · · · · · · · ·			· · · · · · · · · · · · · · · · · · ·											· · · · · · · · · · · · · · · · · · ·								· · · · · · · · · · · · · · · · · · ·		
••••				•••••	••••••		· · · · · · · · · · · · · · · · · · ·											· · · · · · · · · · · · · · · · · · ·				· · · · · · · · · · · · · · · · · · ·			· · · · · · · · · · · · · · · · · · ·		• • • • • • • • •	
· · · · ·	•																											• •
	· · · · · · · · · · · · · · · · · · ·			•	• • • • • • •		· · · · · · · · · · · · · · · · · · ·																		•			
				•			· · · · · · · · · · · · · · · · · · ·																					
· · · · · ·				· · · · · · · · · · · · · · · · · · ·																		· · · · · · · · · · · · · · · · · · ·			· · · · · · · · · · · · · · · · · · ·			
· · · · · · ·							* * * * *												· · · · · · · · · · · · · · · · · · ·			,			· · · · · · · · · · · · · · · · · · ·			
				· · · · · · · · · · · · · · · · · · ·			* * * * * *																					
							· · · · · · · · · · · · · · · · · · ·										1	NORN	IANHU	RST	воч	S' HIGH	SCHC	DOL	· · · · · · · · · · · · · · · · · · ·		• • • • • • • • •	

The line of best fit can be used to make predictions given a value of one of the two variables.

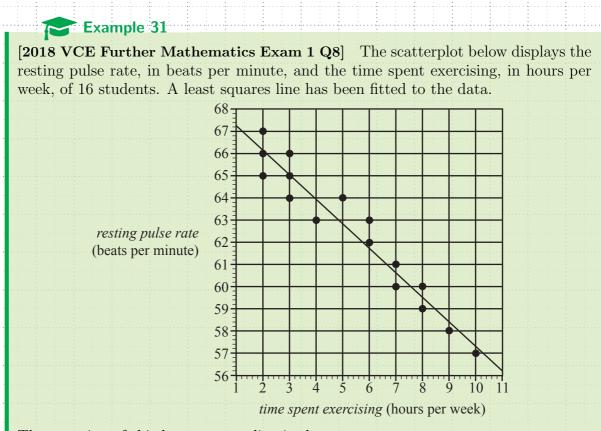
49





- (ii) Use the line of fit to estimate the height difference between a student who wears a size 7.5 shoe and one who wears a size 9 shoe.
- (iii) A student calculated the correlation coefficient to be 1 for this set of data. Explain why this cannot be correct.

Answer: (ii) 11 cm



The equation of this least squares line is closest to

- (A) resting pulse rate = $67.2 0.91 \times$ time spent exercising
- (B) resting pulse rate = $67.2 1.10 \times$ time spent exercising
- (C) resting pulse rate = $68.3 0.91 \times$ time spent exercising
- (D) resting pulse rate = $68.3 1.10 \times$ time spent exercising
- (E) resting pulse rate = $67.2 + 1.10 \times$ time spent exercising

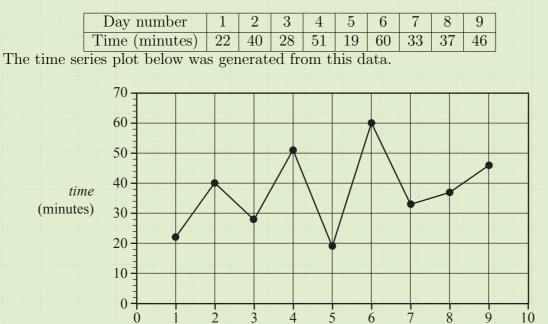
Answer: D

NORMANHURST BOYS' HIGH SCHOOL

51

[2019 VCE Further Mathematics Exam 1 Q14] The time, in minutes, that Liv ran each day was recorded for nine days.

These times are shown in the table below.





A least squares line is to be fitted to the time series plot shown above. The equation of this least squares line, with day number as the explanatory variable, is closest to

(E)

- (A) day number = $23.8 + 2.29 \times \text{time}$ (D)
 - time = $23.8 + 2.29 \times \text{day number}$
- (B) day number = $28.5 + 1.77 \times \text{time}$

time = $28.5 + 1.77 \times$ day number

(C) time = $23.8 + 1.77 \times$ day number

NORMANHURST BOYS' HIGH SCHOOL

52

t

Example	33				
	er Mathematics Exam	n 1 Q13]]	The statistic:	al analysis of a	ı set
of bivariate data ir	nvolving variables x and	y resulted i	in the inform	nation displaye	d in
the table below.					
	Mean	$\overline{x} = 27.8$	$\overline{y} = 33.4$		
	Standard deviation	$\sigma_x = 2.33$	$\sigma_y = 3.24$		
	Equation of the least	y = -2.8	4 + 1.31x		
	squares line				
Using this information data is closest to	tion, the value of the cor	relation coef	ficient r for	this set of bivar	riate
(A) 0.88 (A)	B) 0.89 (C) 0	.92 (E) 0.94	(E) 0.97	

(A	.)	0.	00		- (-	\mathbf{D}	, C	1.0	9		\cup	0.5	2		(D))	0.3	94		(\mathbf{L}))	0.	91	
1.1.1																								
100 A 100																								
and the second second																								
100 A 100																								
100 Aug. 100																								
100 Aug. 100																								
100 Aug. 100																								
and the second second																								

[2019 VCE Further Mathematics Exam 1 Q11] A study was conducted to investigate the effect of drinking coffee on sleep.

In this study, the amount of sleep, in hours, and the amount of coffee drunk, in cups, on a given day were recorded for a group of adults.

The following summary statistics were generated.

	Sleep (hours)	Coffee (cups)	
Mean	7.08	2.42	
Standard deviation	1.12	1.56	
Correlation coefficient (r)	-0.	770	

On average, for each additional cup of coffee drunk, the amount of sleep

 (\mathbf{A}) decreased by 0.55 hours

Example 34

- increased by 1.1 hours (\mathbf{D})
- (B) decreased by 0.77 hours
- (E)increased by 2.3 hours
- (C)decreased by 1.1 hours

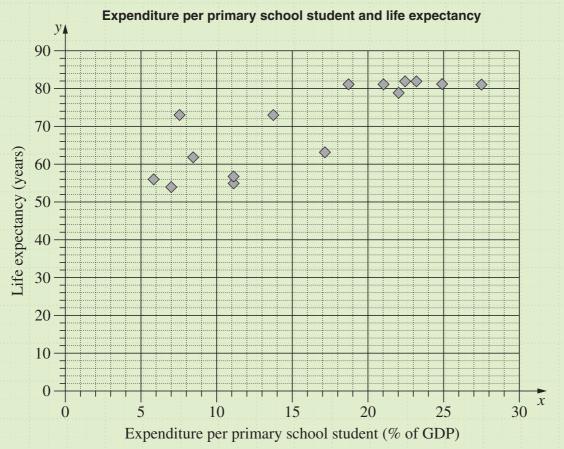
Answer: A

NORMANHURST BOYS' HIGH SCHOOL

53

Answer: D

[2014 General 2 HSC Q30] The scatterplot shows the relationship between expenditure per primary school student, as a percentage of a country's Gross Domestic Product (GDP), and the life expectancy in years for 15 countries.



- i. For the given data, the correlation coefficient, r, is 0.83. What does this indicate about the relationship between expenditure per primary school student and life expectancy for the 15 countries?
- ii. For the data representing expenditure per primary school student, Q_L is 8.4 and Q_U is 22.5.

What is the interquartile range?

iii. Another country has an expenditure per primary school student of 47.6% of its GDP. Would this country be an outlier for this set of data? Justify your answer with calculations.

54

scatterplot are:

[2014 General 2 HSC Q30] Example 35 on the preceding page continued...

iv.

The expenditures per primary school student for the 15 countries in the 2

 $5.9 \quad 7 \quad 7.6 \quad 8.4 \quad 11.2 \quad 11.2 \quad 13.7 \quad 17.1 \quad 18.7$ $21.1 \quad 22 \quad 22.5 \quad 23.2 \quad 24.9 \quad 27.6$

Complete the table below by calculating the mean, \overline{x} , and the standard deviation, σ_x , of these data. Calculate both values to two decimal places. The table also shows the mean, \overline{y} , and the standard deviation, σ_y , of life expectancy for the same 15 countries.

	Mean	Standard deviation
Expenditure per primary	$\overline{x} =$	$\sigma_x =$
school student		
Life expectancy	$\overline{y} = 70.73$	$\sigma_y = 10.94$

v. Using the values from the table in part (iv), show that the equation of the least-squares line of best fit is

$$y = 1.29x + 49.9$$

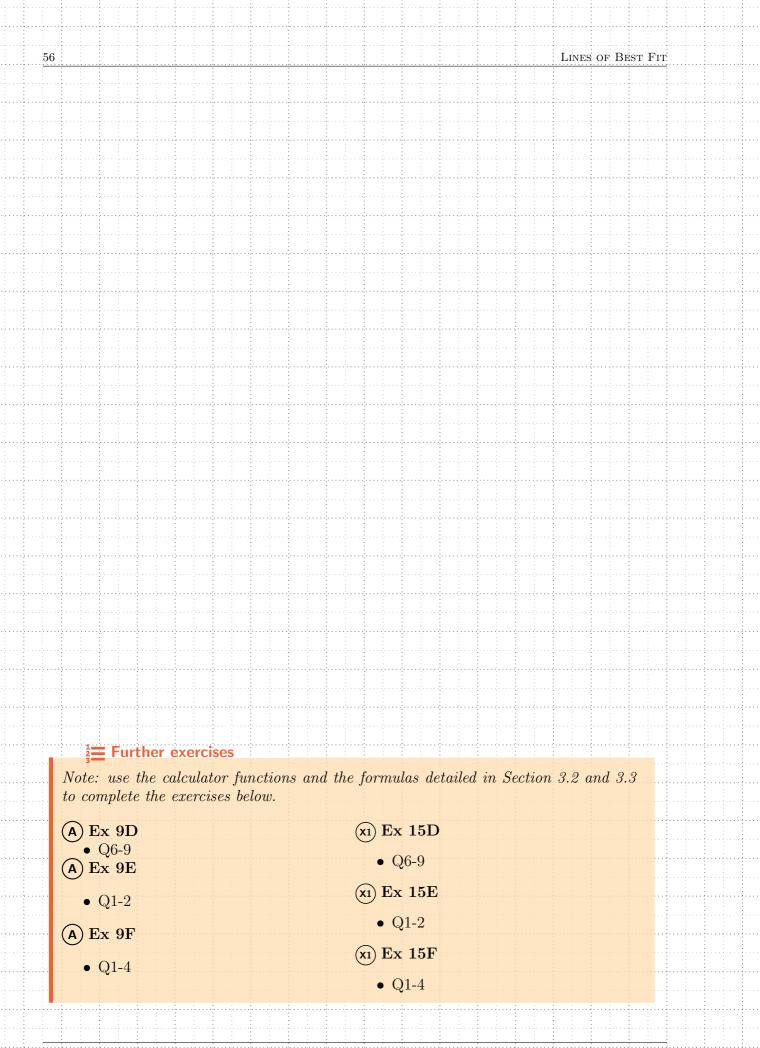
vi.

On the scatterplot provided, draw the least-squares line of best fit,

$$y = 1.29x + 49.9$$

- vii. Using this line, or otherwise, estimate the life expectancy in a country which has an expenditure per primary school student of 18% of its GDP.
- viii. Why is this line NOT useful for predicting life expectancy in a country which has expenditure per primary school student of 60% of its GDP?

Answer: (ii) 14.1 (iii) Yes (iv) $\overline{x} = 16.14$, $\sigma_x = 7.03$ (vi) 73 years



NORM	A'N	HURS	T: B	OYS'	HIG	HISO	CHOO

Section 4

Privacy, Bias and Ethics



Knowledge Understand the effect of bias at different stages of data collection

Important note

🕰 Skills

Identify issues involving bias in real world data, particularly in the media

Vunderstanding

Recognise the ethical issues involved with data collection and analysis

✓ By the end of this section am I able to:
 31.13 Construct, interpret and analyse scatterplots for bivariate numerical data in practical contexts

Digital data collected from a variety of sources involving digital devices, raise *privacy* and *ethical* issues.

The main *ethical issuses* that researchers must consider are:

- **Consent**: Have users agreed to the collection and use of their data?
- **Privacy**: Can the data be used to identify users?
- **Ownership**: Who owns the data and who has the right to determine what the data can be used for?
- Data sharing and reuse: Can the data collected be shared and can the data be used for different purposes than the purpose for which it was originally collected?

Bias can occur throughout the investigative process, from data collection through to statistical analysis and reporting.



Important note

Bias can be in favour or against an outcome, and can skew data and distort the truth of findings. This influences the validity and reliability of conclusions, which in turn impact how decisions are made.

Fill in the spaces

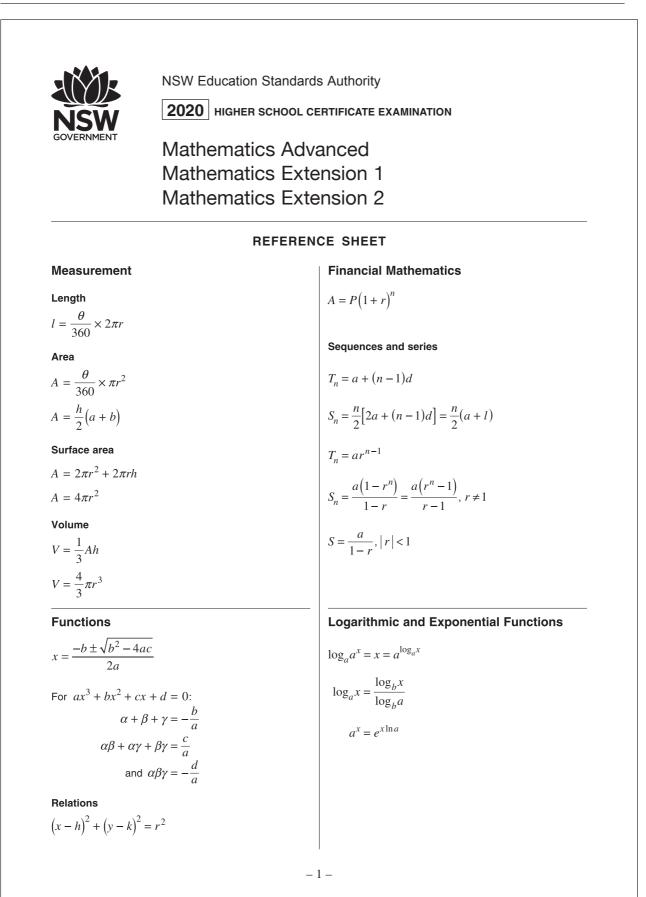
Bias can result from the researcher themselves, as a result of:

- selecting bias <u>sample</u> groups (that do not represent the target population)
- their <u>behaviour</u> (mannerisms, style of dress, speaking tone or body language)
- asking biased <u>questions</u> (phrase questions that influence participants' answers)
- personal opinions (affects their analysis and reporting)

Example 36

) (URL) How coronavirus charts can mislead us

NESA Reference Sheet – calculus based courses



Trigonometric Functions

 $\sin A = \frac{\text{opp}}{\text{hyp}}, \quad \cos A = \frac{\text{adj}}{\text{hyp}}, \quad \tan A = \frac{\text{opp}}{\text{adj}}$ $A = \frac{1}{2}ab\sin C$ $\frac{\sqrt{2}}{\frac{a}{\sin A}} = \frac{b}{\sin B} = \frac{c}{\sin C}$ $\frac{\sqrt{2}}{45^{\circ}}$ $C^{2} = a^{2} + b^{2} - 2ab\cos C$ $\cos C = \frac{a^{2} + b^{2} - c^{2}}{2ab}$ $l = r\theta$ $A = \frac{1}{2}r^{2}\theta$ $\frac{60^{\circ}}{1}$

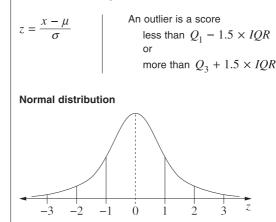
Trigonometric identities

$$\sec A = \frac{1}{\cos A}, \ \cos A \neq 0$$
$$\csc A = \frac{1}{\sin A}, \ \sin A \neq 0$$
$$\cot A = \frac{\cos A}{\sin A}, \ \sin A \neq 0$$
$$\cos^2 x + \sin^2 x = 1$$

Compound angles

 $\sin(A + B) = \sin A \cos B + \cos A \sin B$ $\cos(A + B) = \cos A \cos B - \sin A \sin B$ $\tan(A + B) = \frac{\tan A + \tan B}{1 - \tan A \tan B}$ If $t = \tan \frac{A}{2}$ then $\sin A = \frac{2t}{1 + t^2}$ $\cos A = \frac{1 - t^2}{1 + t^2}$ $\tan A = \frac{2t}{1 - t^2}$ $\cos A \cos B = \frac{1}{2} [\cos(A - B) + \cos(A + B)]$ $\sin A \sin B = \frac{1}{2} [\cos(A - B) - \cos(A + B)]$ $\sin A \cos B = \frac{1}{2} [\sin(A + B) + \sin(A - B)]$ $\cos A \sin B = \frac{1}{2} [\sin(A + B) - \sin(A - B)]$ $\sin^2 nx = \frac{1}{2} (1 - \cos 2nx)$ $\cos^2 nx = \frac{1}{2} (1 + \cos 2nx)$

Statistical Analysis



- approximately 68% of scores have z-scores between -1 and 1
- approximately 95% of scores have z-scores between -2 and 2
- approximately 99.7% of scores have z-scores between –3 and 3

$$E(X) = \mu$$

 $\sqrt{3}$

$$Var(X) = E[(X - \mu)^2] = E(X^2) - \mu^2$$

Probability

$$P(A \cap B) = P(A)P(B)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, P(B) \neq 0$$

Continuous random variables

$$P(X \le x) = \int_{a}^{x} f(x) dx$$
$$P(a < X < b) = \int_{a}^{b} f(x) dx$$

Binomial distribution

$$P(X = r) = {}^{n}C_{r}p^{r}(1-p)^{n-r}$$

$$X \sim \operatorname{Bin}(n, p)$$

$$\Rightarrow P(X = x)$$

$$= {n \choose x}p^{x}(1-p)^{n-x}, x = 0, 1, \dots, n$$

$$E(X) = np$$

$$\operatorname{Var}(X) = np(1-p)$$

- 2 -

Differential Calculus		Integral Calculus
Function	Derivative	$\int f'(x) [f(x)]^n dx = \frac{1}{n+1} [f(x)]^{n+1} + c$
$y = f(x)^n$	$\frac{dy}{dx} = nf'(x)[f(x)]^{n-1}$	$\int \frac{1}{n+1} \frac{n+1}{n+1} \frac{1}{n+1}$ where $n \neq -1$
y = uv	$\frac{dy}{dx} = u\frac{dv}{dx} + v\frac{du}{dx}$	$\int f'(x)\sin f(x)dx = -\cos f(x) + c$
y = g(u) where $u = f(x)$	$\frac{dy}{dx} = \frac{dy}{du} \times \frac{du}{dx}$	$\int f'(x)\cos f(x)dx = \sin f(x) + c$
$y = \frac{u}{v}$	$\frac{dy}{dx} = \frac{v\frac{du}{dx} - u\frac{dv}{dx}}{v^2}$	$\int f'(x)\sec^2 f(x)dx = \tan f(x) + c$
$y = \sin f(x)$	$\frac{dy}{dx} = f'(x)\cos f(x)$	$\int f'(x)e^{f(x)}dx = e^{f(x)} + c$
$y = \cos f(x)$	$\frac{dy}{dx} = -f'(x)\sin f(x)$	
$y = \tan f(x)$	$\frac{dy}{dx} = f'(x)\sec^2 f(x)$	$\int \frac{f'(x)}{f(x)} dx = \ln f(x) + c$
$y = e^{f(x)}$	$\frac{dy}{dx} = f'(x)e^{f(x)}$	$\int f'(x)a^{f(x)}dx = \frac{a^{f(x)}}{\ln a} + c$
$y = \ln f(x)$	$\frac{dy}{dx} = \frac{f'(x)}{f(x)}$	$\int \frac{f'(x)}{\sqrt{a^2 - [f(x)]^2}} dx = \sin^{-1} \frac{f(x)}{a} + c$
$y = a^{f(x)}$	$\frac{dy}{dx} = (\ln a)f'(x)a^{f(x)}$	$\int f'(x) = 1 \int f(x)$
$y = \log_a f(x)$	$\frac{dy}{dx} = \frac{f'(x)}{(\ln a)f(x)}$	$\int \frac{f'(x)}{a^2 + [f(x)]^2} dx = \frac{1}{a} \tan^{-1} \frac{f(x)}{a} + c$
$y = \sin^{-1} f(x)$	$\frac{dy}{dx} = \frac{f'(x)}{\sqrt{1 - \left[f(x)\right]^2}}$	$\int u \frac{dv}{dx} dx = uv - \int v \frac{du}{dx} dx$
$y = \cos^{-1} f(x)$	$\frac{dy}{dx} = -\frac{f'(x)}{\sqrt{1 - [f(x)]^2}}$	$\int_{a}^{b} f(x) dx$
$y = \tan^{-1} f(x)$	$\frac{dy}{dx} = \frac{f'(x)}{1 + [f(x)]^2}$	$\approx \frac{b-a}{2n} \Big\{ f(a) + f(b) + 2 \Big[f(x_1) + \dots + f(x_{n-1}) \Big] \Big\}$ where $a = x_0$ and $b = x_n$
- 3 -		

Combinatorics

$${}^{n}P_{r} = \frac{n!}{(n-r)!}$$

$$\binom{n}{r} = {}^{n}C_{r} = \frac{n!}{r!(n-r)!}$$

$$(x+a)^{n} = x^{n} + \binom{n}{1}x^{n-1}a + \dots + \binom{n}{r}x^{n-r}a^{r} + \dots + a^{n}$$

Vectors

$$\begin{split} \left| \begin{array}{c} \underline{u} \right| &= \left| \begin{array}{c} x\underline{i} + y\underline{j} \right| = \sqrt{x^2 + y^2} \\ \\ \underline{u} \cdot \underline{v} &= \left| \begin{array}{c} \underline{u} \right| \left| \begin{array}{c} \underline{v} \right| \cos \theta = x_1 x_2 + y_1 y_2 \,, \\ \\ \\ \text{where } \begin{array}{c} \underline{u} &= x_1 \underline{i} + y_1 \underline{j} \\ \\ \\ \\ \text{and } \begin{array}{c} \underline{v} &= x_2 \underline{i} + y_2 \underline{j} \end{array} \end{split}$$

 $r_{\tilde{z}} = a + \lambda b_{\tilde{z}}$

Complex Numbers

 $z = a + ib = r(\cos\theta + i\sin\theta)$ $= re^{i\theta}$ $\left[r(\cos\theta + i\sin\theta)\right]^n = r^n(\cos n\theta + i\sin n\theta)$ $= r^n e^{in\theta}$

Mechanics

$$\frac{d^2x}{dt^2} = \frac{dv}{dt} = v\frac{dv}{dx} = \frac{d}{dx}\left(\frac{1}{2}v^2\right)$$
$$x = a\cos(nt + \alpha) + c$$
$$x = a\sin(nt + \alpha) + c$$
$$\ddot{x} = -n^2(x - c)$$

© 2018 NSW Education Standards Authority