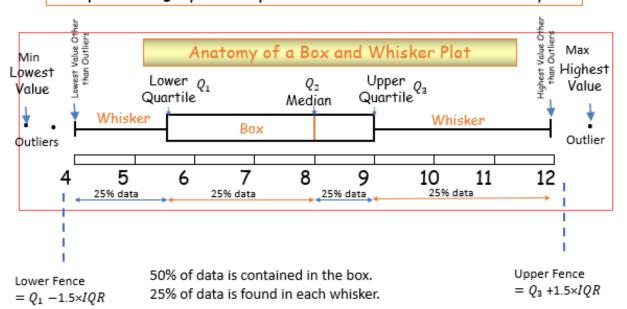


Note: Textbook Summary Notes Section #, Report Instruction Notes #. CAS Instruction Notes #

Box and Whisker Plots

Box plots are graphical representations of 5 number summary.



0.75 ≤ r ≤ 1

Strong, positive, linear association

0.5 ≤ r < 0.75

Moderate, positive, linear association

0.25 ≤ r < 0.5

Weak, positive, linear association

-0.25 < r < 0.25

No association

$-0.5 < r \le -0.25$

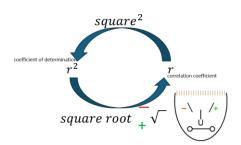
Weak, negative, linear association

-0.75 < r ≤ -0.5

Moderate, negative, linear association

$-1 \le r \le -0.75$

Strong, negative, linear association



Name of Variation	Symbol of Variation	Equation of Variation	CAS	CAS		
			Slope/Equation	Plot/Equation		
Direct Variation	$y \propto x$	y = kx	3.1	y = kx + c		
Inverse Variation	$y \propto \frac{1}{x}$	$y = \frac{k}{x}$	<mark>3.1</mark>	$y = \frac{k}{x} + c$		
Logarithm Variation	$y \propto \log x$	$y = k \log x$	3.1	$y = k \log x + c$		
Variation involving powers	$y \propto x^2$	$y = kx^2$	3.1	$y = kx^2 + c$		
Other variations	$y \propto \sqrt{x}$ $y \propto x^3$ $y \propto \frac{1}{x^2}$	$y = k\sqrt{x}$ $y = kx^{3}$ $y = \frac{k}{x^{2}}$	3.1			

Chapter 9

Data Summary Notes

1A: Types of data

Categorical: characteristics/qualities

Nominal: grouped according to characteristics

Ordinal: can be grouped and ordered

Numerical: numbers/quantities

Discrete: whole numbers, can be counted

Continuous: is measured

1B: displaying categorical data

Count frequency: number of times the category

appears in the data

Percentage frequency: $\frac{count\ frequency}{total\ count} \times 100$

total count

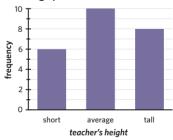
Mode: most frequently occurring value or category

Frequency Table:

teacher's	frequency					
height	number	%				
short	6	25.0				
average	10	41.7				
tall	8	33.3				
total	24	100.0				

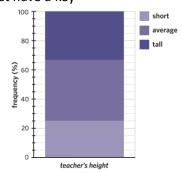
Bar chart:

Must have gaps between bars



Segmented bar chart:

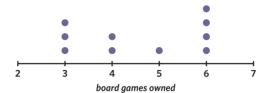
- Can be count or percentage frequency
- Must have a key



1C: Displaying Numerical data

Dot plot

- Discrete data
- Small data sets



Stem and leaf plot

- Needs a key
- Can have class intervals (splitting the stem in two if it is really large)

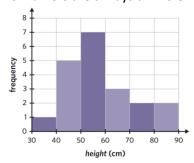
	two in tells really large,													
No intervals			With class intervals of 5											
Key: 1 2 = 1.2			Key: 1 2 = 12											
1	3	3	4	6	8			0	1	1	2	3	4	
2	0	4	9					0	5	6	6	8	8	9
3	1	1	1	4	5	8		1	2	3	3			
4	2							1	6	7	7	8	9	9

Grouped frequency tables

height (cm)	frequency					
neight (Cili)	number	%				
30-<40	1	5				
40-<50	5	25				
50-<60	7	35				
60-<70	3	15				
70-<80	2	10				
80-<90	2	10				
total	20	100				

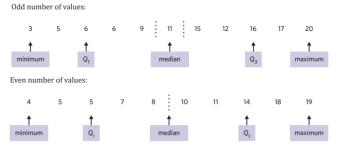
Histogram:

- Continuous data
- Intervals no gaps between bars
- No gaps between bars
- X-axis markers are always a whole number



1E: the five-number summary and boxplots 5 number summary:

- Minimum: smallest value in data set
- Q1: median of the lower half
- Median: middle value in an ordered data set
- Q3: median of the upper half
- Maximum: largest value in data set



Spread: refers to how variable the data set it Range = maximum - minimum

Interquartile range: measure of spread of the middle 50% of a data set. Accurate measure of spread when outliers are present.

$$IQR = Q_3 - Q_1$$

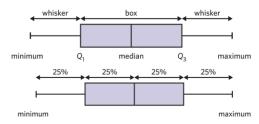
Outliers: values which fall outside of what is 'normal'. Outliers are still the minimum and maximum value! Fence: defines the boundary of what is an outlier. If a value is less than the lower fence or greater than the upper fence it is considered to be an outlier.

$$lower fence = Q_1 - (1.5 \times IQR)$$

$$upper fence = Q_3 + (1.5 \times IQR)$$

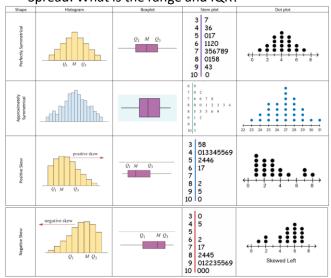
$$1.5 \times IQR$$

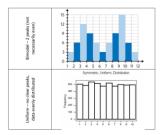
Boxplots:



1F: describing numerical data

- Shape: is the data symmetrical, skewed or have any outliers?
- Centre: What is the median value?
- Spread: What is the range and IQR?



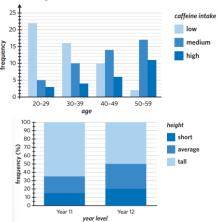


2A: association between 2 variables

Two-way frequency table:

- Columns = EV, Rows = RV
- Percentage frequency is used for greater accuracy when making comparisons if sample sizes are different

Grouped and segmented bar charts:



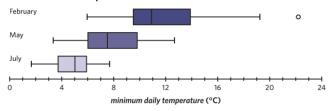
Describing the association between two variables:

- Whether or not an association between the two variables exists
- Appropriate percentages to support findings

2B: association between numerical and categorical variables

• Back to back stem plot

Parallel boxplot



 Making comparisons: refer to 1F and compare shape, centre and spread of the two categories

2C: association between two numerical variables

Response variable: RV, may be explained or predicted by changes in the explanatory variable.

Explanatory variable: EV, used to explain or predict the changes observed in the response variable.

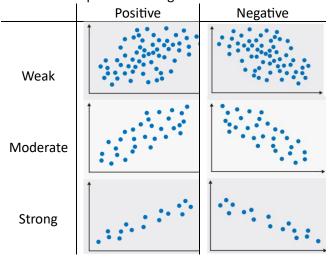
• 'EV explains the RV'

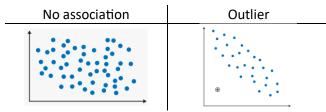
Scatterplots:

• EV = x axis, RV = y axis

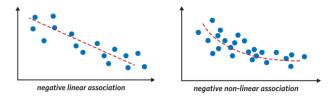
Describing relationship/analysing scatterplots:

- Strength: how close the data points are together
- Direction: positive or negative





• Form: linear (straight) or non-linear (curved)



2D: Correlation and causation

Pearson's correlation coefficient (r): numerical value that determines strength and direction between two numerical variables, assuming:

- Data is linear
- Data is numeric
- No outliers present

$0.75 \le r \le 1$	Strong, positive, linear association
$0.5 \le r < 0.75$	Moderate, positive, linear association
$0.25 \le r < 0.5$	Weak, positive, linear association
-0.25 < r < 0.25	No association
$-0.5 < r \le -0.25$	Weak, negative, linear association
$-0.75 < r \le -0.5$	Moderate, negative, linear association
$-1 \le r \le -0.75$	Strong, negative, linear association

3A: fitting a least squares regression line

Least squares regression line (LSRL): is the line which creates the minimum sum of the squares of residuals. There are assumptions:

Data is numerical

- The relationship between variables is linear
- There are no clear outliers present

The line is used to show the general trend in the data and is given by the equation:

y = a + bxIntercept Slope

Determining LSRL from a graph: Find the intercept (a) and the slope (b).

- Intercept: read directly from the graph when the EV is 0
- Slope: choose two points on the line that you can clearly ready the coordinates. Use the rule:

$$b = \frac{rise}{run} = \frac{y_2 - y_1}{x_2 - x_1}$$

Drawing the LSRL on a graph: Sub in the first value on the x-axis and the last value on the x-axis into the equation. Plot the two points, join the line using a ruler.

3B: Interpreting LSRL: use the following statements, fill in EV and RV and values of a and b.

y-intercept: when the EV is 0, the RV is a.

Slope: for every one-unit increase in the EV, the RV increases/decreases by b. (If b is positive, increases, if b is negative, decreases)

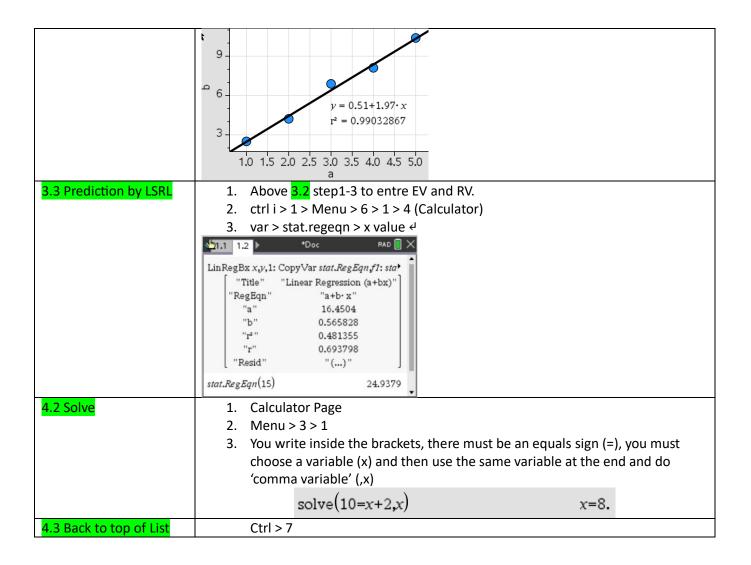
Making predictions: the LSRL can be used to predict the value of the RV from the EV.

Interpolation: predicting within the range of data. **Extrapolation:** predicting outside the range of data. Less reliable.

How to predict: sub the EV value that you are predicting for into the LSRL equation to predict the RV.

CAS reference sheet: Data

1.1 Univariate Data	1 C+rl	n > 1 Ento	r data in flict a	and enroadehoe	ats' langura a titla	\			
Statistics: to find the 5	 Ctrl n > 4, Enter data in 'list and spreadsheets' (ensure a title) Menu > 4 > 1 > 1 OR ctrl i > 1 > menu > 6 > 1 > 1 in a new page 								
number summary,				data in the 'X-L		,			
mean and standard	I diaz					Variable Reference ▶			
deviation	MinX	1.	+			3.			
	QıX	1.5	+		Σχ	15.			
	MedianX	3.							
	Q ₃ X	4.5			Σx²	55.			
	MaxX	5.			SX := Sn	1 . 58113			
1.2 Graphing Univariate	1. Ctrl	n > 4, Entei	data in 'list a	nd spreadshee	ts' (ensure a title)				
Data: dot plot, boxplots,			i <i>OR</i> home m	enu > 5 to ope	n a new page and	select 'data			
histograms		statistics'							
			ta to the x-axis	s pe you are afte	\r				
	Histogram s		ise the plot ty	pe you are arte	:1				
	_	nu > 2 > 2 >	2 > 1						
				interval size. A	lignment is where	e the graph will			
	star	t on the x-a	xis		Equal Bin Width Settings	-			
		m: menu >	_		Alignment 0				
	8. Fred	quency or P	ercent: menu	> 2 > 2 > 1	OK Cano	cet			
2.1 Graphing Bivariate		ate first Bar							
Data: Grouped Bar Chart	Menu >	2 > 9 > sele	ect second set	of data					
2.2 Graphing Bivariate	2. Crea	 2. Create first boxplot 3. Menu > 2 > 5 > select second set of data 							
Data: parallel boxplots	-								
3.1 Bivariate Data Statistics: a, b, r, r ²	1. Ctrl a tit		r data in 'list a	nd spreadshee	ts'(ensure you giv	e both columns			
	2. Mei	nu > 4 > 1 >	4 (Spreadshee	ets) <i>OR</i> ctrl i > :	1 > Menu > 6 > 1 >	> 4 (Calculator)			
	3. Sele	ct your EV i	in the 'X-List' a	and your RV in	the 'Y-List'				
	Linear Regression	on (a+hv)							
		ist: 'a		RegEqn	a+b*x				
		ist: 'b	<u> </u>	а	0.51				
	Save RegEqr	to: f1	▶	L	1.07				
	Frequency I		>	b	1.97				
	Category I			r²	0.99032				
			OK Cancel	r	0.99515				
3.2 Scatterplots		-		nd spreadshee	ts' (ensure you giv	ve both			
		mns a title							
	2. Ctrl doc OR ctrl i OR home menu > 5 to open a new page and select 'data								
	and statistics' 3. Click to add the EV to the x-axis and the RV to the y-axis								
	3. Click to add the LV to the A-dals and the IV to the y-dals								
	Scatterplot Least Squares Regression Line:								
	4. Menu 4 > 6 > 2 <i>OR</i> ctrl i > 1 > Menu > 6 > 1 > 4 (Calculator)								



1a. Univariate Categorical Data: Frequency Table

The [types of categories] of [total frequency] [frequency type] were classified as [list of categories].

Modal Category

The majority of [frequency type], [modal percentage], were found to be [modal category].

Of the remaining [frequency types], [percentage X] were found to be [category X], while [percentage Y] were found to be [category Y], and while [etc.].

Equal Categories

The [frequency types] all had roughly the same percentages where [category X] had [percentage X], [category Y] had [percentage Y], [etc.].

1b. Univariate Numerical Data: Histogram, Dot Plot, Stem Plot

The shape of the distribution is [symmetric/positively skewed/negatively skewed]

Refer to 1F: Describing numerical data

The distribution has a [standard dev./range/IQR] of [value]

The distribution has a [mean/median/mode] of [value]

The distribution [has #/has no] outliers.

1c. Univariate Numerical Data: Box Plot

The distribution is [positively skewed/negatively skewed] with [outliers/no outliers]. The distribution is centred at [value], the median value. The spread of the distribution, is measured by the IQR, is [value] and, as measured by the range [value]. If outliers present: There are [value] many outliers: [list of outliers]

2a. Bivariate Data (Both Categorical): Two-way Frequency Table

Worked example: Is there an association between interest in sports and age group? Yes, the percentage of males with a high level of interest in sport steadily decreases with age group from 56.5 % for the 'under 18 years' age group, to 35.0% for the '36-50 years' age group.



2b. Bivariate Data (one categorical, One Numerical): Comparing two boxplots:

The distributions at [variable name] are [symmetric/positively skewed/negatively skewed] for both [boxplot variables]. There [are/are no] outliers. The median [variable name] is higher for [boxplot 1], (M= value), than [boxplot 2], (M= value). The IQR is also greater for [boxplot 1], (IQR= value), than [boxplot 2], (IQR= value). The range of [variable name] is also greater for [boxplot 1], (R= value), than [boxplot 2], (R= value). Yes, there is an association between [variable name1] [variable name 2] as their median change from (M= value) to (M= value).

2c. Bivariate Data (Both Numerical): Scatter Plot

There is a [strong/moderate/weak], [positive/negative], [linear/non-linear] relationship between [response variable y] and [explanatory variable x]. There [are/are no] clear outliers.

2d. The coefficient of determination (r²):

The coefficient of determination indicates that $[r^2 \times 100]$ % of the variation in $[response \ variable]$ is explained by the variation in $[explanatory \ variable]$ and [remaining %] is explained by other factors and not explained by $[explanatory \ variable]$.

2e. Least squares line:

The equation of the regression line is: [response variable] = [a] + [b] x [explanatory variable]

Slope (b):

On average, [response variable] [increases/decreases] by [b units] for every one [unit] increase in [explanatory variable].

y- intercept (a):

When [explanatory variable] is 0, [response variable] is predicted to be [a units].

2f. Residual Plot

The residual plot shows a [random scatter/ curved pattern] indicating there is a [linear/non-linear] relationship between [response variable] and [explanatory variable].

2g. Prediction Reliability

The prediction using [explanatory variable] of [value of EV] is [reliable/ unreliable] as it is [within/ outside] the date range and is therefore [interpolation/extrapolation].