

5. Correlation and Regression

Mr Faruk

Teacher of Mathematics
BSc/MSc/PGCE Mathematics

✉ cieigcsolutions@gmail.com



4. Sara was studying the relationship between rainfall, r mm, and humidity, h %, in the UK. She takes a random sample of 11 days from May 1987 for Leuchars from the large data set.

She obtained the following results.

h	93	86	95	97	86	94	97	97	87	97	86
r	1.1	0.3	3.7	20.6	0	0	2.4	1.1	0.1	0.9	0.1

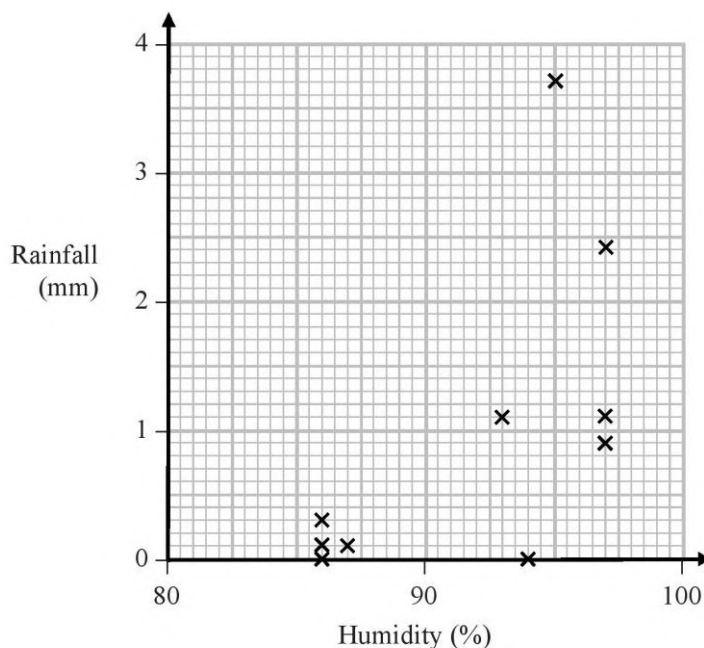
Sara examined the rainfall figures and found

$$Q_1 = 0.1 \quad Q_2 = 0.9 \quad Q_3 = 2.4$$

A value that is more than 1.5 times the interquartile range (IQR) above Q_3 is called an outlier.

- (a) Show that $r = 20.6$ is an outlier. (1)
- (b) Give a reason why Sara might (i) include
(ii) exclude
 this day's reading. (2)

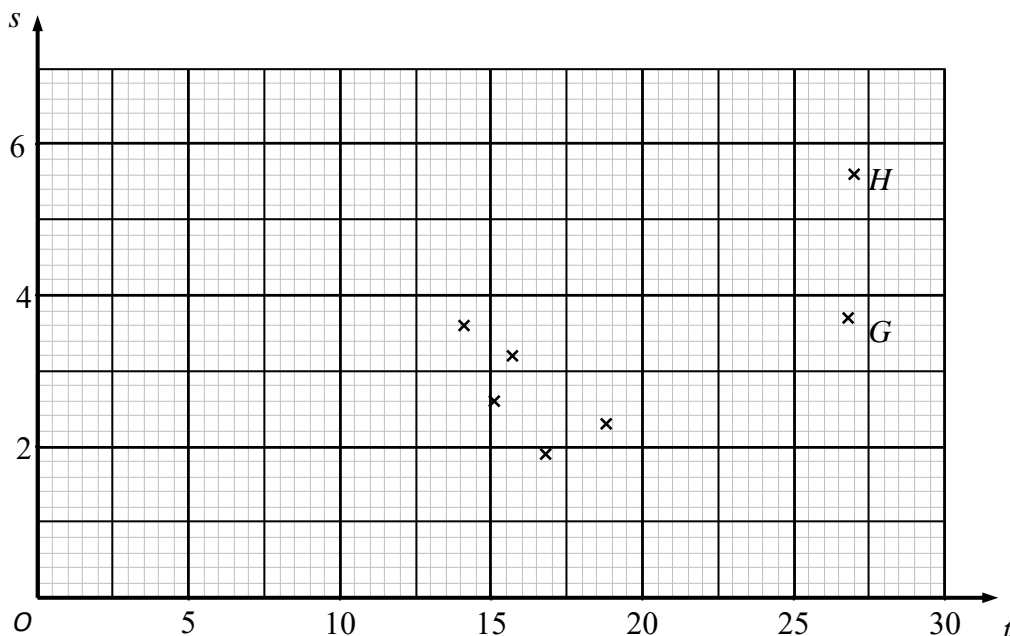
Sara decided to exclude this day's reading and drew the following scatter diagram for the remaining 10 days' values of r and h .



- (c) Give an interpretation of the correlation between rainfall and humidity. (1)

Specimen 2018 AL Statistics

2. A researcher believes that there is a linear relationship between daily mean temperature and daily total rainfall. The 7 places in the northern hemisphere from the large data set are used. The mean of the daily mean temperatures, t °C, and the mean of the daily total rainfall, s mm, for the month of July in 2015 are shown on the scatter diagram below.



- (a) With reference to the scatter diagram, explain why a linear regression model may not be suitable for the relationship between t and s . (1)

The researcher calculated the product moment correlation coefficient for the 7 places and obtained $r = 0.658$.

- (b) Stating your hypotheses clearly, test at the 10% level of significance, whether or not the product moment correlation coefficient for the population is greater than zero. (3)
- (c) Using your knowledge of the large data set, suggest the names of the 2 places labelled G and H . (1)
- (d) Using your knowledge from the large data set, and with reference to the locations of the two places labelled G and H , give a reason why these places have the highest temperatures in July. (2)
- (e) Suggest how you could make better use of the large data set to investigate the relationship between daily mean temperature and daily total rainfall. (1)

(Total 7 marks)

1. A company is introducing a job evaluation scheme. Points (x) will be awarded to each job based on the qualifications and skills needed and the level of responsibility. Pay (£ y) will then be allocated to each job according to the number of points awarded.

Before the scheme is introduced, a random sample of 8 employees was taken and the linear regression equation of pay on points was $y = 4.5x - 47$

- (a) Describe the correlation between points and pay. (1)
- (b) Give an interpretation of the gradient of this regression line. (1)
- (c) Explain why this model might not be appropriate for all jobs in the company. (1)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

3. Barbara is investigating the relationship between average income (GDP per capita), x US dollars, and average annual carbon dioxide (CO_2) emissions, y tonnes, for different countries.

She takes a random sample of 24 countries and finds the product moment correlation coefficient between average annual CO_2 emissions and average income to be 0.446

- (a) Stating your hypotheses clearly, test, at the 5% level of significance, whether or not the product moment correlation coefficient for all countries is greater than zero.

(3)

Barbara believes that a non-linear model would be a better fit to the data.

She codes the data using the coding $m = \log_{10}x$ and $c = \log_{10}y$ and obtains the model $c = -1.82 + 0.89m$

The product moment correlation coefficient between c and m is found to be 0.882

- (b) Explain how this value supports Barbara's belief.

(1)

- (c) Show that the relationship between y and x can be written in the form $y = ax^n$ where a and n are constants to be found.

(5)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

2. Jerry is studying visibility for Camborne using the large data set June 1987.

The table below contains two extracts from the large data set.

It shows the daily maximum relative humidity and the daily mean visibility.

Date	Daily Maximum Relative Humidity	Daily Mean Visibility
Units	%	
10/06/1987	90	5300
28/06/1987	100	0

(The units for Daily Mean Visibility are deliberately omitted.)

Given that daily mean visibility is given to the nearest 100,

- (a) write down the range of distances in metres that corresponds to the recorded value 0 for the daily mean visibility.

(1)

Jerry drew the following scatter diagram, Figure 2, and calculated some statistics using the June 1987 data for Camborne from the large data set.

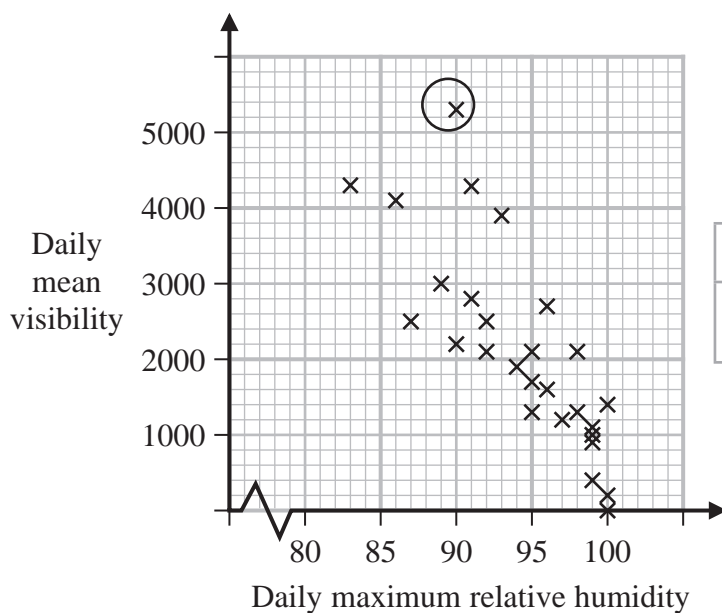


Figure 2

	Q_1	IQR
Daily mean visibility	1100	1600
Daily maximum relative humidity (%)	92	8

Jerry defines an outlier as a value that is more than 1.5 times the interquartile range above Q_3 or more than 1.5 times the interquartile range below Q_1 .

- (b) Show that the point circled on the scatter diagram is an outlier for visibility.

(2)

- (c) Interpret the correlation between the daily mean visibility and the daily maximum relative humidity.

(1)

DO NOT WRITE IN THIS AREA

Jerry drew the following scatter diagram, Figure 3, using the June 1987 data for Camborne from the large data set, but forgot to label the x -axis.

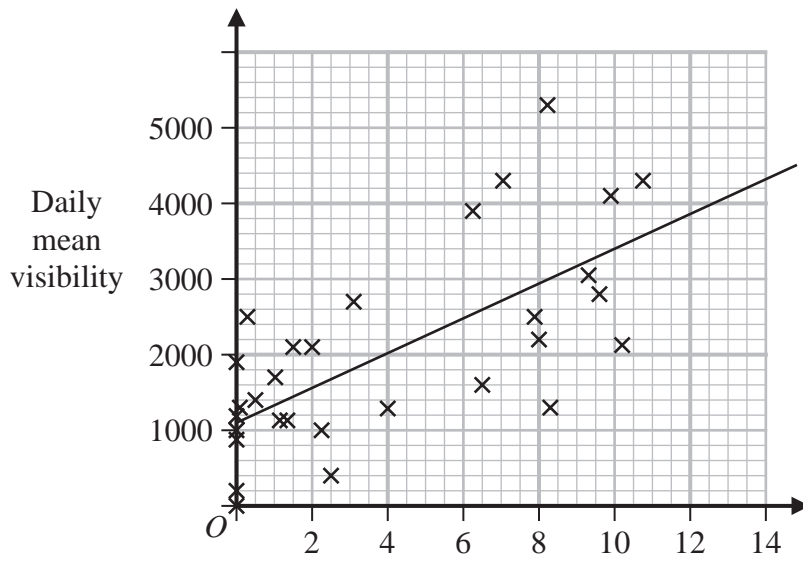


Figure 3

(d) Using your knowledge of the large data set, suggest which variable the x -axis on this scatter diagram represents.

(1)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

2. A random sample of 15 days is taken from the large data set for Perth in June and July 1987. The scatter diagram in Figure 1 displays the values of two of the variables for these 15 days.

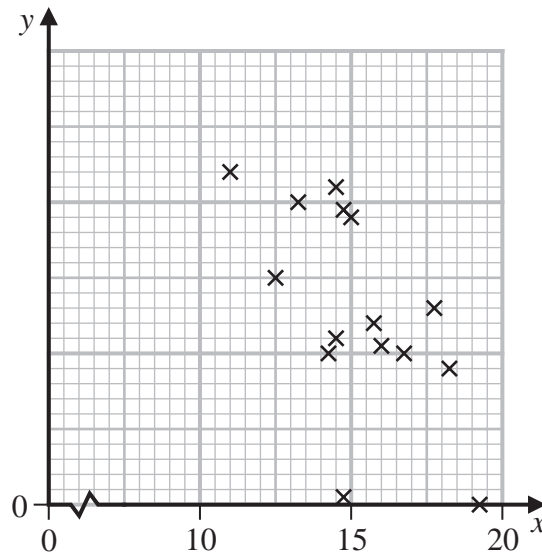


Figure 1

- (a) Describe the correlation.

(1)

The variable on the x -axis is Daily Mean Temperature measured in $^{\circ}\text{C}$.

- (b) Using your knowledge of the large data set,

- suggest which variable is on the y -axis,
- state the units that are used in the large data set for this variable.

(2)

Stav believes that there is a correlation between Daily Total Sunshine and Daily Maximum Relative Humidity at Heathrow.

He calculates the product moment correlation coefficient between these two variables for a random sample of 30 days and obtains $r = -0.377$

- (c) Carry out a suitable test to investigate Stav's belief at a 5% level of significance. State clearly

- your hypotheses
- your critical value

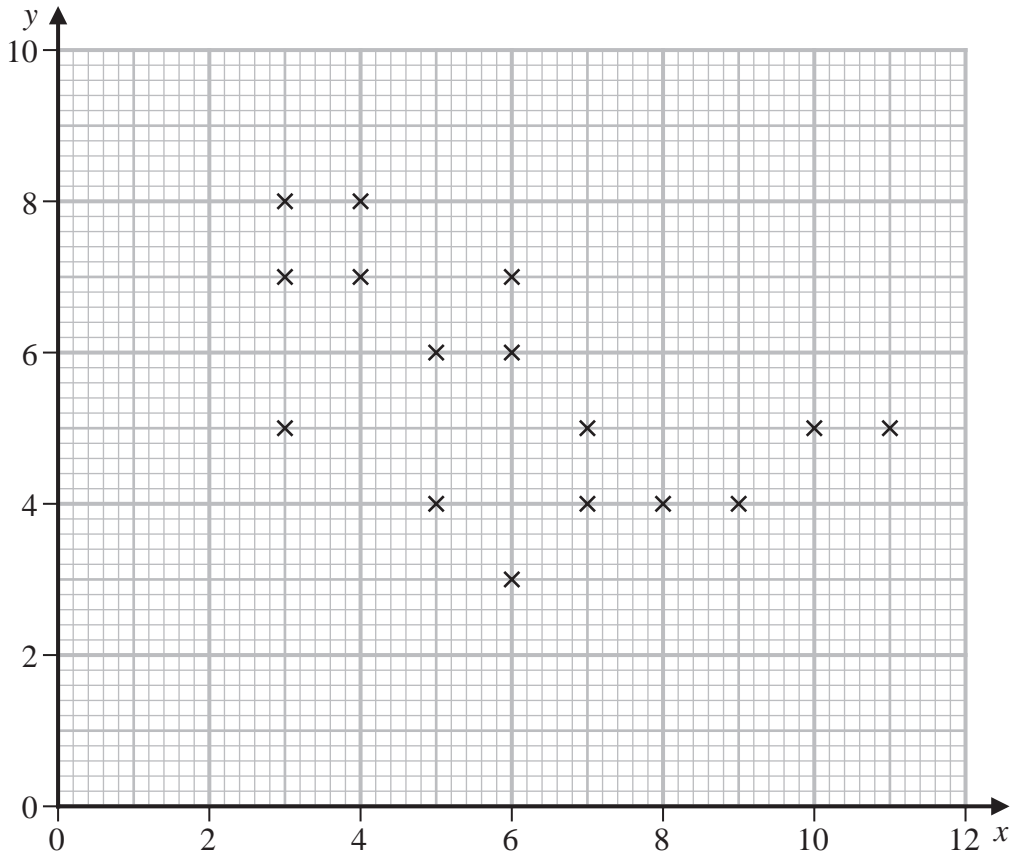
(3)

On a random day at Heathrow the Daily Maximum Relative Humidity was 97%

- (d) Comment on the number of hours of sunshine you would expect on that day, giving a reason for your answer.

(1)

Question 2 continued.



DO NOT WRITE IN THIS AREA

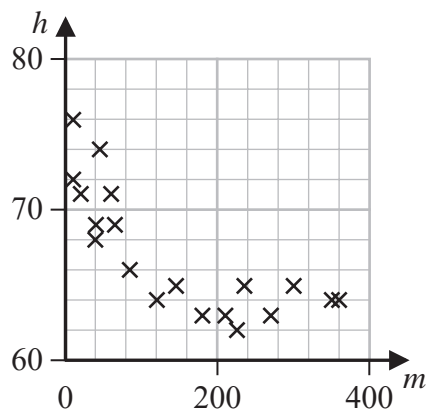
DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

6. Anna is investigating the relationship between exercise and resting heart rate. She takes a random sample of 19 people in her year at school and records for each person

- their resting heart rate, h beats per minute
- the number of minutes, m , spent exercising each week

Her results are shown on the scatter diagram.



(a) Interpret the nature of the relationship between h and m

(1)

Anna codes the data using the formulae

$$x = \log_{10} m$$

$$y = \log_{10} h$$

The product moment correlation coefficient between x and y is -0.897

(b) Test whether or not there is significant evidence of a negative correlation between x and y

You should

- state your hypotheses clearly
- use a 5% level of significance
- state the critical value used

(3)

The equation of the line of best fit of y on x is

$$y = -0.05x + 1.92$$

(c) Use the equation of the line of best fit of y on x to find a model for h on m in the form

$$h = am^k$$

where a and k are constants to be found.

(5)

2. Fred and Nadine are investigating whether there is a linear relationship between Daily Mean Pressure, p hPa, and Daily Mean Air Temperature, t °C, in Beijing using the 2015 data from the large data set.

Fred randomly selects one month from the data set and draws the scatter diagram in Figure 1 using the data from that month.

The scale has been left off the horizontal axis.

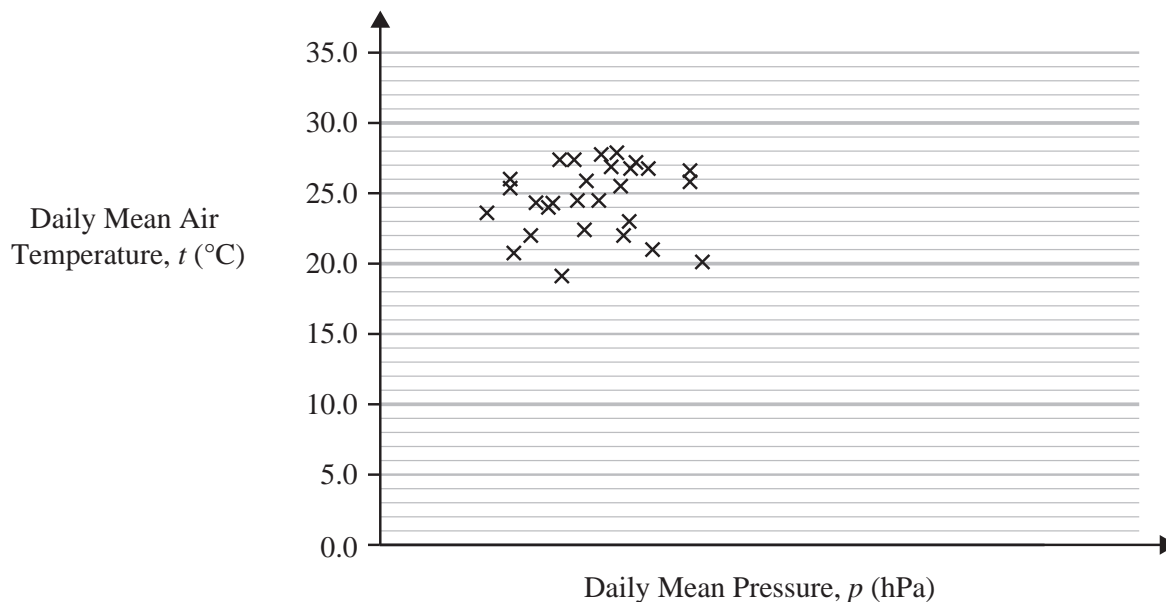


Figure 1

- (a) Describe the correlation shown in Figure 1.

(1)

Nadine chooses to use all of the data for Beijing from 2015 and draws the scatter diagram in Figure 2.

She uses the same scales as Fred.

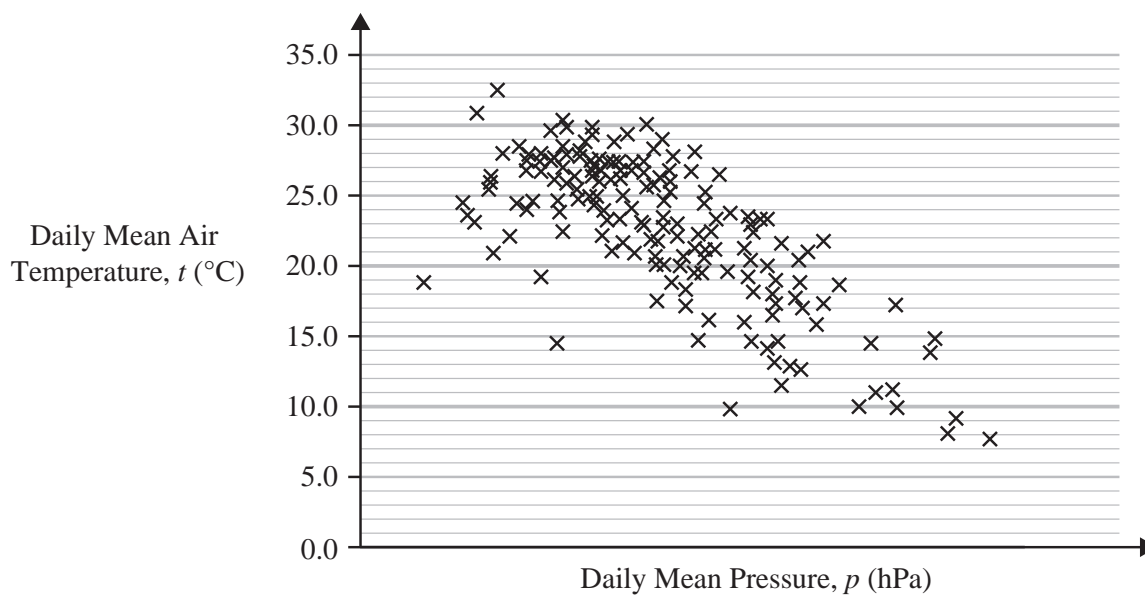


Figure 2

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

2. Amar is studying the flight of a bird from its nest.

He measures the bird's height above the ground, h metres, at time t seconds for 10 values of t

Amar finds the equation of the regression line for the data to be $h = 38.6 - 1.28t$

- (a) Interpret the gradient of this line.

(1)

The product moment correlation coefficient between h and t is -0.510

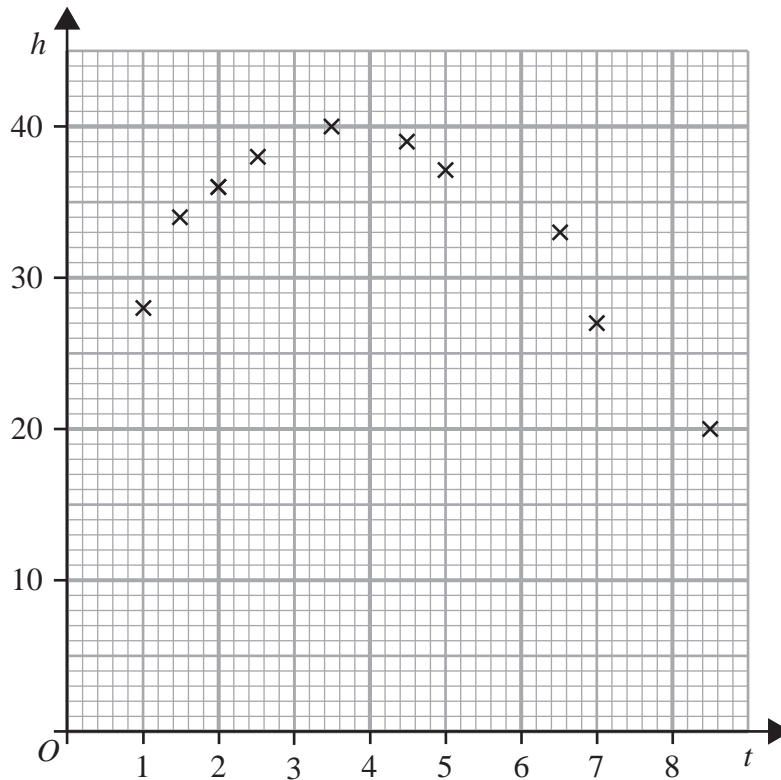
- (b) Test whether or not there is evidence of a negative correlation between the height above the ground and the time during the flight.

You should

- state your hypotheses clearly
- use a 5% level of significance
- state the critical value used

(3)

Jane draws the following scatter diagram for Amar's data.



- (c) With reference to the scatter diagram, state, giving a reason, whether or not the regression line $h = 38.6 - 1.28t$ is an appropriate model for these data.

(1)

Jane suggests an improved model using the variable $u = (t - k)^2$ where k is a constant.

She obtains the equation $h = 38.1 - 0.78u$

- (d) Choose a suitable value for k to write Jane's improved model for h in terms of t only.

(1)

