

2014-2018

Chapter 5:

Correlation and

Regression

Mr Faruk

Teacher of Mathematics
BSc/MSc/PGCE Mathematics

✉ cieigcsolutions@gmail.com



3. Jean works for an insurance company. She randomly selects 8 people and records the price of their car insurance, £ p , and the time, t years, since they passed their driving test. The data is shown in the table below.

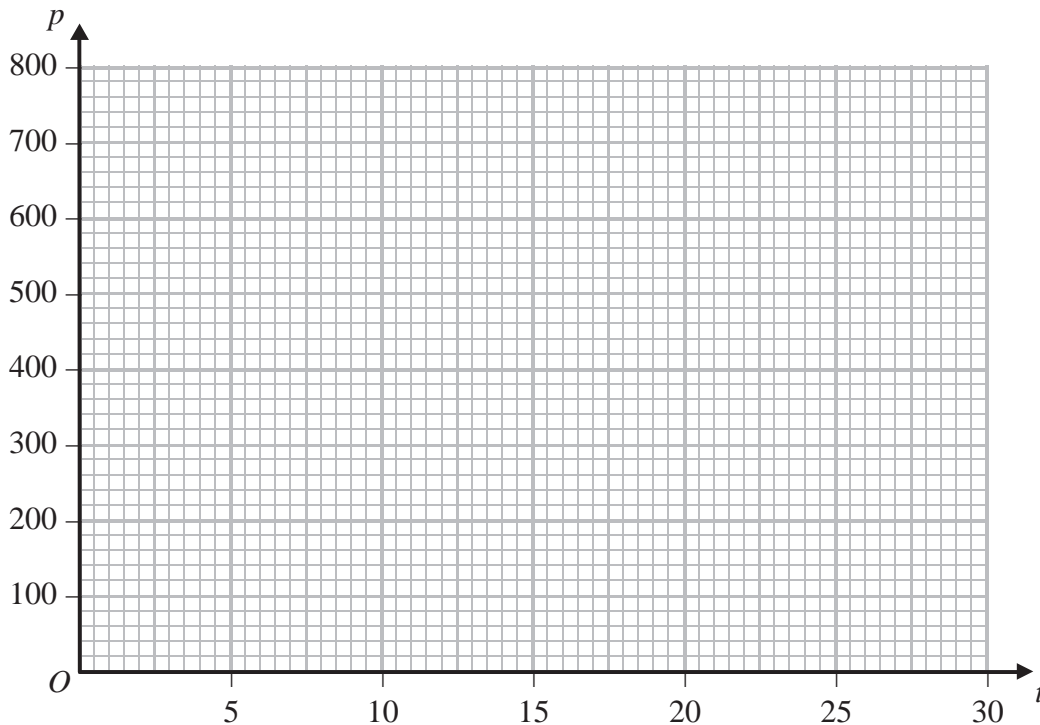
t	10	13	17	18	22	24	25	27
p	720	650	430	490	500	390	280	300

(You may use $\bar{t} = 19.5$, $\bar{p} = 470$, $S_{tp} = -6080$, $S_{tt} = 254$, $S_{pp} = 169\,200$)

- (a) On the graph below draw a scatter diagram for these data. (2)
- (b) Comment on the relationship between p and t . (1)
- (c) Find the equation of the regression line of p on t . (4)
- (d) Use your regression equation to estimate the price of car insurance for someone who passed their driving test 20 years ago. (2)

Jack passed his test 39 years ago and decides to use Jean's data to predict the price of his car insurance.

- (e) Comment on Jack's decision. Give a reason for your answer. (2)



3. The table shows the price of a bottle of milk, m pence, and the price of a loaf of bread, b pence, for 8 different years.

m	29	29	35	39	41	43	44	46
b	75	83	91	121	120	126	119	126

(You may use $S_{bb} = 3083.875$ and $S_{mm} = 305.5$)

- (a) Find the exact value of $\sum bm$ (1)
- (b) Find S_{bm} (3)
- (c) Calculate the product moment correlation coefficient between b and m (2)
- (d) Interpret the value of the correlation coefficient. (1)

A ninth year is added to the data set. In this year the price of the bottle of milk is 46 pence and the price of a loaf of bread is 175 pence.

- (e) Without further calculation, state whether the value of the product moment correlation coefficient will increase, decrease or stay the same when all nine years are used. Give a reason for your answer. (2)

3. A publisher collects information about the amount spent on advertising, £ x , and the sales, y books, for some of her publications. She collects information for a random sample of 8 textbooks and codes the data using $v = \frac{x + 50}{200}$ and $s = \frac{y}{1000}$ to give

v	0.60	8.10	4.30	0.40	1.60	6.40	2.50	5.10
s	1.84	6.73	5.95	1.30	2.45	7.46	4.82	6.25

[You may use: $\sum v = 29$ $\sum s = 36.8$ $\sum s^2 = 209.72$ $\sum vs = 177.311$ $S_{vv} = 55.275$]

- (a) Find S_{vs} and S_{ss} (3)
- (b) Calculate the product moment correlation coefficient for these data. (2)

The publisher believes that a linear regression model may be appropriate to describe these data.

- (c) State, giving a reason, whether or not your answer to part (b) supports the publisher's belief. (1)
- (d) Find the equation of the regression line of s on v , giving your answer in the form $s = a + bv$ (4)
- (e) Hence find the equation of the regression line of y on x for the sample of textbooks, giving your answer in the form $y = c + dx$ (3)

The publisher calculated the regression line for a sample of novels and obtained the equation

$$y = 3100 + 1.2x$$

She wants to increase the sales of books by spending more money on advertising.

- (f) State, giving your reasons, whether the publisher should spend more money on advertising textbooks or novels. (2)

3. A scientist measured the salinity of water, x g/kg, and recorded the temperature at which the water froze, y °C, for 12 different water samples. The summary statistics are listed below.

$$\sum x = 504 \quad \sum y = -27 \quad \sum x^2 = 22842 \quad \sum y^2 = 62.98$$

$$\sum xy = -1190.7 \quad S_{xx} = 1674 \quad S_{yy} = 2.23$$

- (a) Find the mean and variance of the recorded temperatures. (3)

Priya believes that the higher the salinity of water, the higher the temperature at which the water freezes.

- (b) (i) Calculate the product moment correlation coefficient between x and y
 (ii) State, with a reason, whether or not this value supports Priya's belief. (4)

- (c) Find the least squares regression line of y on x in the form $y = a + bx$
 Give the value of a and the value of b to 3 significant figures. (4)

- (d) Estimate the temperature at which water freezes when the salinity is 32 g/kg (1)

The coding $w = 1.8y + 32$ is used to convert the recorded temperatures from °C to °F

- (e) Find an equation of the least squares regression line of w on x in the form $w = c + dx$ (2)
- (f) Find (3)
- (i) the variance of the recorded temperatures when converted to °F
- (ii) the product moment correlation coefficient between w and x

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

5. Franca is the manager of an accountancy firm. She is investigating the relationship between the salary, £ x , and the length of commute, y minutes, for employees at the firm. She collected this information from 9 randomly selected employees.

The salary of each employee was then coded using $w = \frac{x - 20\,000}{1000}$

The table shows the values of w and y for the 9 employees.

w	6	8	8	-1	25	15	3	-2	19
y	45	50	35	65	25	40	50	75	20

(You may use $\sum w = 81$ $\sum y = 405$ $\sum wy = 2490$ $S_{ww} = 660$ $S_{yy} = 2500$)

- (a) Calculate the salary of the employee with $w = -2$ (1)
- (b) Show that, to 3 significant figures, the value of the product moment correlation coefficient between w and y is -0.899 (3)
- (c) State, giving a reason, the value of the product moment correlation coefficient between x and y (1)

The least squares regression line of y on w is $y = 60.75 - 1.75w$

- (d) Find the equation of the least squares regression line of y on x giving your answer in the form $y = a + bx$ (3)
- (e) Estimate the length of commute for an employee with a salary of £21 000 (2)

Franca uses the regression line to estimate the length of commute for employees with salaries between £25 000 and £40 000

- (f) State, giving a reason, whether or not these estimates are reliable. (2)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

2. Paul believes there is a relationship between the value and the floor size of a house. He takes a random sample of 20 houses and records the value, £ v , and the floor size, s m²

The data were coded using $x = \frac{s - 50}{10}$ and $y = \frac{v}{100\,000}$ and the following statistics obtained.

$$\sum x = 441.5, \quad \sum y = 59.8, \quad \sum x^2 = 11\,261.25, \quad \sum y^2 = 196.66, \quad \sum xy = 1474.1$$

- (a) Find the value of S_{xy} and the value of S_{xx} (3)

- (b) Find the equation of the least squares regression line of y on x in the form $y = a + bx$ (3)

The least squares regression line of v on s is $v = c + ds$

- (c) Show that $d = 1020$ to 3 significant figures and find the value of c (3)

- (d) Estimate the value of a house of floor size 130 m² (2)

- (e) Interpret the value d (1)

Paul wants to increase the value of his house. He decides to add an extension to increase the floor size by 31 m²

- (f) Estimate the increase in the value of Paul's house after adding the extension. (1)

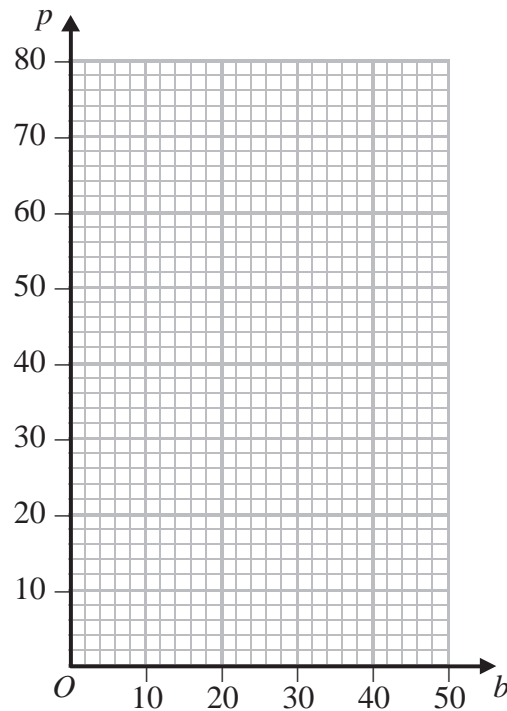
7. A doctor is investigating the correlation between blood protein, p , and body mass index, b .

He takes a random sample of 8 patients and the data are shown in the table below.

Patient	A	B	C	D	E	F	G	H
b	32	36	40	44	42	21	27	37
p	18	21	31	39	21	12	19	70

(a) Draw a scatter diagram of these data on the axes provided.

(2)



The doctor decides to leave out patient H from his calculations.

(b) Give a reason for the doctor's decision.

(1)

For the 7 patients A, B, C, D, E, F and G ,

$$S_{bp} = 369, \quad S_{pp} = 490 \quad \text{and} \quad S_{bb} = 423\frac{5}{7}$$

(c) Find the product moment correlation coefficient, r , for these 7 patients.

(2)

(d) Without any further calculations, state how r would differ from your answer in part (c) if it was calculated for all 8 patients.

(1)

5. Tomas is studying the relationship between temperature and hours of sunshine in *Seapron*. He records the midday temperature, t °C, and the hours of sunshine, s hours, for a random sample of 9 days in October. He calculated the following statistics

$$\sum s = 15 \quad \sum s^2 = 44.22 \quad \sum t = 127 \quad S_{tt} = 10.89$$

- (a) Calculate S_{ss} (2)

Tomas calculated the product moment correlation coefficient between s and t to be 0.832 correct to 3 decimal places.

- (b) State, giving a reason, whether or not this correlation coefficient supports the use of a linear regression model to describe the relationship between midday temperature and hours of sunshine. (1)

- (c) State, giving a reason, why the hours of sunshine would be the explanatory variable in a linear regression model between midday temperature and hours of sunshine. (1)

- (d) Find S_{st} (3)

- (e) Calculate a suitable linear regression equation to model the relationship between midday temperature and hours of sunshine. (4)

- (f) Calculate the standard deviation of s (1)

Tomas uses this model to estimate the midday temperature in *Seapron* for a day in October with 5 hours of sunshine.

- (g) State the value of Tomas' estimate. (1)

Given that the values of s are all within 2 standard deviations of the mean,

- (h) comment, giving your reason, on the reliability of this estimate. (2)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

1. A random sample of 10 cars of different makes and sizes is taken and the published miles per gallon, p , and the actual miles per gallon, m , are recorded. The data are coded using variables $x = \frac{p}{10}$ and $y = m - 25$

The results for the coded data are summarised below.

x	6.89	3.67	5.92	5.04	4.87	3.92	4.71	5.14	3.65	5.23
y	30	3	22	15	13	8	15	13.5	3	19

(You may use $\sum y^2 = 2628.25$ $\sum xy = 768.58$ $S_{xx} = 9.25924$ $S_{.y} = 74.664$)

- (a) Show that $S_{yy} = 626.025$ (2)
- (b) Find the product moment correlation coefficient between x and y . (2)
- (c) Give a reason to support fitting a regression model of the form $y = a + bx$ to these data. (1)
- (d) Find the equation of the regression line of y on x , giving your answer in the form $y = a + bx$.
Give the value of a and the value of b to 3 significant figures. (3)

A car's published miles per gallon is 44

- (e) Estimate the actual miles per gallon for this particular car. (3)
- (f) Comment on the reliability of your estimate in part (e). Give a reason for your answer. (2)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

4. A doctor is studying the scans of 30-week old foetuses. She takes a random sample of 8 scans and measures the length, f mm, of the leg bone called the femur. She obtains the following results.

52 53 56 57 57 59 60 62

- (a) Show that $S_{ff} = 80$ (3)

The doctor also measures the head circumference, h mm, of each foetus and her results are summarised as

$$\sum h = 2209 \quad \sum h^2 = 610\,463 \quad S_{fh} = 182$$

- (b) Find S_{hh} (2)

- (c) Calculate the product moment correlation coefficient between the length of the femur and the head circumference for these data. (2)

The doctor believes that there is a linear relationship between the length of the femur and the head circumference of 30-week old foetuses.

- (d) State, giving a reason, whether or not your calculation in part (c) supports the doctor's belief. (1)

- (e) Find an equation of the regression line of h on f . (4)

The doctor plans in future to measure the femur length, f , and then use the regression line to estimate the corresponding head circumference, h .

A statistician points out that there will always be the chance of an error between the true head circumference and the estimated value of the head circumference.

Given that the error, E mm, has the normal distribution $N(0, 4^2)$

- (f) find the probability that the estimate is within 3 mm of the true value. (3)

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

DO NOT WRITE IN THIS AREA

