

## Linear regression Mixed exercise

$$1 \text{ a } S_{st} = \sum st - \frac{\sum s \sum t}{n} = 31\,185 - \frac{553 \times 549}{12} = 31\,185 - 25\,299.75 = 5885.25$$

$$b = \frac{S_{st}}{S_{ss}} = \frac{5885.25}{6193} = 0.95030\dots = 0.950 \text{ (3 s.f.)}$$

$$a = \bar{t} - b\bar{s} = 45.75 - (0.95030\dots \times 46.0833) = 1.95672\dots = 1.96 \text{ (3 s.f.)}$$

Hence equation of regression line of  $t$  on  $s$  is:  $t = 1.96 + 0.95s$

$$b \text{ } t = 1.9567\dots + (0.9503\dots \times 50) = 49.4717 = 49.5 \text{ (3 s.f.)}$$

$$2 \text{ a } S_{xx} = \sum x^2 - \frac{(\sum x)^2}{n} = 97.73 - \frac{33.1 \times 33.1}{12} = 6.4291\dots = 6.429 \text{ (4 s.f.)}$$

$$S_{xy} = \sum xy - \frac{\sum x \sum y}{n} = 195.94 - \frac{33.1 \times 66.8}{12} = 195.94 - 184.26 = 11.68333\dots = 11.68 \text{ (4 s.f.)}$$

$$b \text{ } \bar{x} = \frac{\sum x}{n} = \frac{33.1}{12} = 2.7583\dots \quad \bar{y} = \frac{\sum y}{n} = \frac{66.8}{12} = 5.5666\dots$$

$$b = \frac{S_{xy}}{S_{xx}} = \frac{11.6833\dots}{6.4291\dots} = 1.8172\dots = 1.82 \text{ (3 s.f.)}$$

$$a = \bar{y} - b\bar{x} = 5.5666\dots - (1.8172\dots \times 2.7583\dots) = 0.55421\dots = 0.554 \text{ (3 s.f.)}$$

Equation is:  $y = 0.554 + 1.82x$

$$c \text{ Length of leaf} = 0.5542\dots + (1.8172\dots \times 3) = 6.0058\dots = 6.01 \text{ cm (3 s.f.)}$$

3 a Calculating the summary statistics gives:

$$\sum x^2 = 43\,622.85 \quad \sum x = 467.1 \quad \sum y = 7805 \quad \sum xy = 666\,045$$

$$S_{xx} = \sum x^2 - \frac{(\sum x)^2}{n} = 43\,622.85 - \frac{467.1 \times 467.1}{8} = 16\,350.048\dots = 16\,350 \text{ (5 s.f.)}$$

$$S_{xy} = 666\,045 - \frac{467.1 \times 7805}{8} = 210\,330.56\dots = 210\,331 \text{ (6 s.f.)}$$

$$b \text{ } \bar{x} = \frac{\sum x}{n} = \frac{467.1}{8} = 58.3875 \quad \bar{y} = \frac{\sum y}{n} = \frac{7805}{8} = 975.625$$

$$b = \frac{S_{xy}}{S_{xx}} = \frac{210\,330.56}{16\,350.048} = 12.8642\dots = 12.86 \text{ (4 s.f.)}$$

$$a = \bar{y} - b\bar{x} = 975.625 - (12.8642\dots \times 58.3875) = 224.5155\dots = 224.5 \text{ (4 s.f.)}$$

Equation is:  $y = 224.5 + 12.86x$

$$c \text{ Gross National Product} = 224.515\dots + (12.8642\dots \times 100) = 1510.93\dots = 1511 \text{ (4 s.f.)}$$

**3 d**  $3500 = 224.515\dots + 12.864\dots x$

$$\Rightarrow \text{Energy consumption } (x) = \frac{3500 - 224.515\dots}{12.8642\dots} = 255 \text{ (3 s.f.)}$$

**e** This answer is likely to be unreliable as it involves extrapolation. The value of 3500 is well outside the limits of the data set used.

**4 a**  $S_{xy} = \sum xy - \frac{\sum x \sum y}{n} = 84.25 - \frac{25.5 \times 13.5}{6} = 84.25 - 57.375 = 26.875$

$$\bar{x} = \frac{\sum x}{n} = \frac{25.5}{6} = 4.25 \quad \bar{y} = \frac{\sum y}{n} = \frac{13.5}{6} = 2.25$$

$$b = \frac{S_{xy}}{S_{xx}} = \frac{26.875}{59.88} = 0.44881\dots = 0.449 \text{ (3 s.f.)}$$

$$a = \bar{y} - b\bar{x} = 2.25 - (0.44881\dots \times 4.25) = 0.3425\dots = 0.343 \text{ (3 s.f.)}$$

Equation is:  $y = 0.343 + 0.449x$

**b**  $t - 2 = 0.3425\dots + 0.4488\dots \left(\frac{m}{2}\right)$

$$\Rightarrow t = 2.3425\dots + 0.2244\dots m$$

$$\Rightarrow t = 2.34 + 0.224m \quad (\text{rounding the parameters to 3 s.f.})$$

**c** Tail length =  $2.3425\dots + (0.2244\dots \times 10) = 4.5865\dots = 4.6 \text{ cm (2 s.f.)}$

**5 a**  $S_{xx} = \sum x^2 - \frac{(\sum x)^2}{n} = 747 - \frac{81 \times 81}{10} = 747 - 656.1 = 90.9$

$$S_{xy} = \sum xy - \frac{\sum x \sum y}{n} = 1413 - \frac{81 \times 151}{10} = 1413 - 1223.1 = 189.9 = 190 \text{ (3 s.f.)}$$

**b**  $b = \frac{S_{xy}}{S_{xx}} = \frac{189.9}{90.9} = 2.0891\dots = 2.09 \text{ (3 s.f.)}$

$$\bar{x} = \frac{\sum x}{n} = \frac{81}{10} = 8.1 \quad \bar{y} = \frac{\sum y}{n} = \frac{151}{10} = 15.1$$

$$a = \bar{y} - b\bar{x} = 15.1 - (2.0891\dots \times 8.1) = -1.8217\dots = -1.82 \text{ (3 s.f.)}$$

Equation is:  $y = -1.82 + 2.09x$

**c**  $\frac{p-50}{2} = -1.8217\dots + 2.0891\dots \left(\frac{r-10}{2}\right)$

$$\Rightarrow p - 50 = 3.6434\dots + 2.0891\dots r - 20.891$$

$$\Rightarrow p = 25.465\dots + 2.0891\dots r$$

$$\Rightarrow p = 25.5 + 2.09r \quad (\text{rounding parameters to 3 s.f.})$$

**d** Pulse rate =  $25.5 + (2.09 \times 22) = 71 \text{ beats per minute (2 s.f.)}$

**5 e** The value of 22 breaths per minute equates to  $x = 6$  and is within range of the data set used. So the answer to part **d** is reasonably reliable since it involves interpolation.

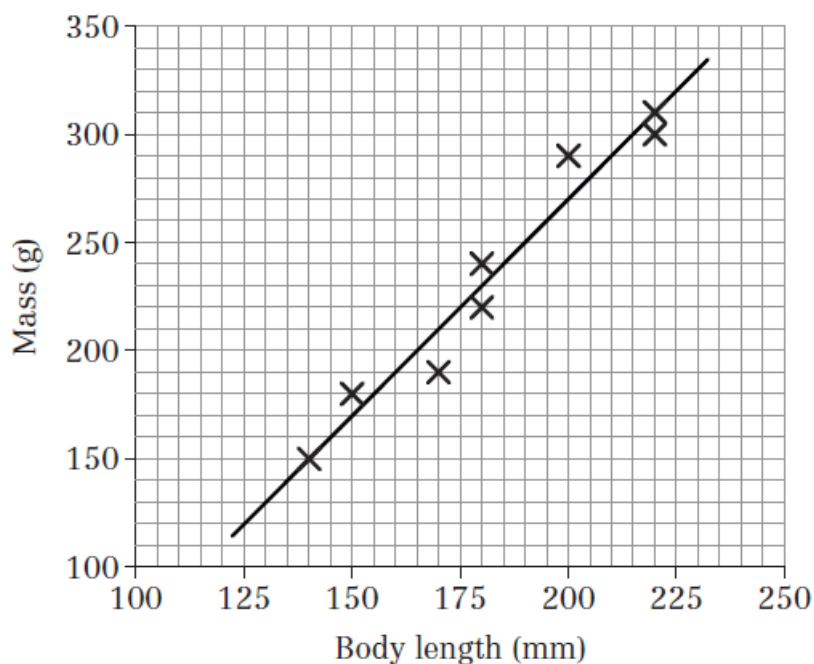
**6 a** The figure of 0.79 is the average amount of food consumed (in kg) in 1 week by 1 hen.

**b**  $y = 0.16 + 0.79 \times 30 = 23.86 = 23.9 \text{ kg}$  (3 s.f.)

**c** Food needed  $= 0.16 + 0.79 \times 50 = 39.66 \text{ kg}$

$$\text{Cost of feed} = \frac{39.66}{10} \times 12 = \text{£}47.592 = \text{£}47.59$$

**7 a** This is a scatter diagram of the data. (The diagram also shows the regression line, found in part **e**.)



**b** There appears to be a linear relationship between body length and body mass.

7 c Calculating the summary statistics for  $l$  and  $w$  gives:

$l$	14	15	17	18	18	20	22	22
$w$	15	18	19	22	24	29	30	31

$$\sum l^2 = 2726 \quad \sum l = 146 \quad \sum w = 188 \quad \sum lw = 3553$$

$$\bar{l} = \frac{\sum l}{n} = \frac{146}{8} = 18.25 \quad \bar{w} = \frac{\sum w}{n} = \frac{188}{8} = 23.5$$

$$S_{ll} = \sum l^2 - \frac{(\sum l)^2}{n} = 2726 - \frac{146 \times 146}{8} = 2726 - 2664.5 = 61.5$$

$$S_{lw} = \sum lw - \frac{\sum l \sum w}{n} = 3553 - \frac{146 \times 188}{8} = 3553 - 3431 = 122$$

$$b = \frac{S_{lw}}{S_{ll}} = \frac{122}{61.5} = 1.9837\dots = 1.98 \text{ (3 s.f.)}$$

$$a = \bar{w} - b\bar{l} = 23.5 - (1.9837\dots \times 18.25) = 23.5 - 36.2032\dots = -12.7032\dots = -12.7 \text{ (3 s.f.)}$$

Equation is:  $w = -12.7 + 1.98l$

d  $\frac{y}{10} = -12.7 + \left(1.98 \times \frac{x}{10}\right) \Rightarrow y = -127 + 1.98x$  (multiply through by 10)

e See diagram for part a.

f Mass =  $-127.0\dots + 1.983\dots \times 210 = 289.43\dots = 290$  grams (2 s.f.)

This is reliable since it involves interpolation. The mass of 210 is within the range of the data.

g Voles B and C are both underweight so were probably removed from the river. Vole A is slightly overweight so was probably left in the river.

8 a Calculating the summary statistics for  $x$  and  $y$  gives:

$x$	0	3	12	5	14	6	9
$y$	7	9	15	9	13	11	13

$$\sum x = 49 \quad \sum x^2 = 491 \quad \sum y = 77 \quad \sum xy = 617$$

$$S_{xy} = \sum xy - \frac{\sum x \sum y}{n} = 617 - \frac{49 \times 77}{7} = 617 - 539 = 78$$

$$S_{xx} = \sum x^2 - \frac{(\sum x)^2}{n} = 491 - \frac{49^2}{7} = 491 - 343 = 148$$

$$8 \text{ b} \quad \bar{x} = \frac{\sum x}{7} = \frac{49}{7} = 7 \quad \bar{y} = \frac{\sum y}{7} = \frac{77}{7} = 11$$

$$b = \frac{S_{xy}}{S_{xx}} = \frac{78}{148} = 0.52702\dots = 0.5270 \text{ (4 s.f.)}$$

$$a = \bar{y} - b\bar{x} = 11 - (0.52702\dots \times 7) = 7.3108\dots = 7.311 \text{ (4 s.f.)}$$

$$\text{Equation is: } y = 7.311 + 0.5270x \quad (\text{parameters to 4 s.f.})$$

$$c \quad \frac{w}{400} = 7.3108\dots + 0.52702\dots(n-10) \text{ (multiply by 400)}$$

$$\Rightarrow w = 816.24\dots + 210.808n$$

$$\Rightarrow w = 816.2 + 210.8n \quad (\text{parameters to 4 s.f.})$$

$$d \quad w = 816.24\dots + 210.808\dots \times 20 = 5032 \text{ kg}$$

e This is far outside the range of values. This is extrapolation.

f This table sets out the residuals for each coded data point ( $x$ ,  $y$ ):

$x$	$y$	$y = 7.311 + 0.527x$	$\epsilon$
0	7	7.311	-0.311
3	9	8.892	0.108
12	15	13.635	1.365
5	9	9.946	-0.946
14	13	14.689	-1.689
6	11	10.473	0.527
9	13	12.054	0.946

g This table sets out the residuals for each coded data point ( $n$ ,  $w$ ):

$n$	$w$	$w = 816.2 + 210.8n$	$\epsilon$
10	2800	2924.2	-124.2
13	3600	3556.6	43.4
22	6000	5453.8	546.2
15	3600	3978.2	-378.2
24	5200	5875.4	-675.4
16	4400	4189	211
19	5200	4821.4	378.6

h They are related by the same code as that used for  $y$  in terms of  $w$ .

9 a This table sets out the residuals for each data point:

$t$	$T$	$T = 80.445 - 4.289t$	$\varepsilon$
2	72	71.867	0.133
3	68	67.578	0.422
5	59	59	0
6	54	54.711	-0.711
7	50	50.422	-0.422
9	42	41.844	0.156
10	38	37.555	0.445

b No: the residuals are not randomly scattered about zero.

$$\text{c } \text{RSS} = S_{TT} - \frac{(S_{Tt})^2}{S_{tt}} = 957.43 - \frac{(-223)^2}{52} = 1.10 \text{ (3 s.f.)}$$

d The second reaction is most likely to have a linear fit since the RSS is lower.

$$10 \text{ a } S_{ss} = \sum s^2 - \frac{(\sum s)^2}{n} = 395.76 - \frac{53.4^2}{8} = 39.315$$

$$S_{sf} = \sum sf - \frac{\sum s \sum f}{n} = 171.66 - \frac{53.4 \times 29.9}{8} = -27.9225$$

$$\text{b } b = \frac{S_{sf}}{S_{ss}} = \frac{-27.9225}{39.315} = -0.71022\dots = -0.710 \text{ (3 s.f.)}$$

$$\bar{s} = \frac{\sum s}{n} = \frac{53.4}{8} = 6.675 \qquad \bar{f} = \frac{\sum f}{n} = \frac{29.9}{8} = 3.7375$$

$$a = \bar{f} - b\bar{s} = 3.7375 - \frac{(-27.9225)}{39.315} \times 6.675 = 8.4782\dots = 8.48 \text{ (3 s.f.)}$$

Hence the equation of the regression line of  $f$  on  $s$  is:  $f = 8.48 - 0.710s$

$$\text{c } f = 8.48 - 0.710 \times 7.5 = 3.2 \text{ mm (2 s.f.)}$$

$$\text{d } S_{ff} = \sum f^2 - \frac{(\sum f)^2}{n} = 131.93 - \frac{29.9^2}{8} = 20.17875$$

$$\text{RSS} = S_{ff} - \frac{(S_{sf})^2}{S_{ss}} = 20.17875 - \frac{(-27.9225)^2}{39.315} = 0.347 \text{ (3 s.f.)}$$

$$\text{e } \sum \varepsilon = 0 \Rightarrow -0.177 - 0.196 + 0.256 + 0.124 + x - 0.129 - 0.216 - 0.032 = 0 \Rightarrow x = 0.37$$

f A linear regression model is not suitable: the residuals go negative, positive, negative so they are not randomly scattered about zero.

$$11 \text{ a } b = \frac{S_{xy}}{S_{xx}} = \frac{75}{60} = 1.25$$

$$\bar{x} = \frac{2+3+4+5+6+7+8+9+10}{9} = 6 \qquad \bar{y} = \frac{4+5+7+8+9+11+12+11+15}{9} = \frac{82}{9}$$

$$a = \bar{y} - b\bar{x} = \frac{82}{9} - 1.25 \times 6 = 1.6111\dots = 1.61 \text{ (3 s.f.)}$$

Hence the equation of the regression line of  $y$  on  $x$  is:  $y = 1.61 + 1.25x$

$$11 \text{ b } \text{RSS} = S_{yy} - \frac{(S_{xy})^2}{S_{xx}} = 98.89 - \frac{75^2}{60} = 5.14$$

11 c This table sets out the residuals for each data point:

$x$	$y$	$y = 1.61 + 1.25x$	$\varepsilon$
2	4	4.11	-0.11
3	5	5.36	-0.36
4	7	6.61	0.39
5	8	7.86	0.14
6	9	9.11	-0.11
7	11	10.36	0.64
8	12	11.61	0.39
9	11	12.86	-1.86
10	15	14.11	0.89

11 d The outlier is the wombat that is nine days old, i.e. data point (9, 11).

11 e i There is no correct answer.

It could be a recording error – so ignore this data point.

It could be a valid data point; the wombat could have been a ‘runt’ or underweight due to illness, for example – so include this reading.

11 e ii Removing the data point gives the following summary statistics:

$$S_{xx} = 49.875 \quad S_{xy} = 68.625 \quad \bar{x} = \frac{45}{8} = 5.625 \quad \bar{y} = \frac{71}{8} = 8.875$$

$$b = \frac{S_{xy}}{S_{xx}} = \frac{68.625}{49.875} = 1.37593\dots = 1.38 \text{ (3 s.f.)}$$

$$a = \bar{y} - b\bar{x} = 8.875 - \frac{68.625}{49.875} \times 5.625 = 1.13533\dots = 1.14 \text{ (3 s.f.)}$$

So the equation of the regression line is now:  $y = 1.14 + 1.38x$

11 e iii Mass of wombat =  $1.135\dots + 1.376\dots \times 20 = 28.65\dots = 28.7 \text{ g (3 s.f.)}$

11 e iv The estimate is unreliable since 20 days lies (well) outside the range of the data.

$$12 \text{ a } S_{tt} = \sum t^2 - \frac{(\sum t)^2}{n} = 42.33 - \frac{17.7^2}{8} = 3.16875$$

$$S_{ts} = \sum ts - \frac{\sum t \sum s}{n} = 42.16 - \frac{17.7 \times 17.5}{8} = 3.44125$$

$$b = \frac{S_{ts}}{S_{tt}} = \frac{3.44125}{3.16875} = 1.0859\dots = 1.09 \text{ (3 s.f.)}$$

$$\bar{t} = \frac{\sum t}{n} = \frac{17.7}{8} = 2.2125 \quad \bar{s} = \frac{\sum s}{n} = \frac{17.5}{8} = 2.1875$$

$$a = \bar{s} - b\bar{t} = 2.1875 - \frac{3.44125}{3.16875} \times 2.2125 = -0.21526\dots = -0.215 \text{ (3 s.f.)}$$

Hence the equation of the regression line of  $s$  on  $t$  is:  $s = -0.215 + 1.09t$

**b** Predicted number of employees ( $s$ ) =  $(-0.215 + 1.09 \times 2.3) \times 100 = 229$  (to nearest whole number)

**c**  $\sum \varepsilon = 0 \Rightarrow 0.0121 - 0.0137 - 0.0395 + 0.1347 - 0.0997 + p + 0.0745 + 0.0487 = 0 \Rightarrow p = -0.1171$   
So  $p = 0.117$  (3 s.f.)

**d** Linear model is suitable since the residuals are randomly scattered about zero.

$$\text{e } S_{ss} = \sum s^2 - \frac{(\sum s)^2}{n} = 42.07 - \frac{17.5^2}{8} = 3.78875$$

$$\text{RSS} = S_{ss} - \frac{(S_{st})^2}{S_{tt}} = 3.78875 - \frac{(3.44125)^2}{3.16875} = 0.0516 \text{ (3 s.f.)}$$

**f** The UK sample is likely to have a better linear fit since the RSS is much smaller.

### Challenge

**a**  $\varepsilon_i = y_i - (a + bx_i)$  where  $i = 1, 2, \dots, n$

$$\sum \varepsilon = \sum (y_i - a - bx_i)$$

Given that  $a = \bar{y} - b\bar{x}$ , this can be written as  $\sum \varepsilon = \sum (y_i - \bar{y} + b\bar{x} - bx_i)$

$$\sum (y_i - \bar{y} + b\bar{x} - bx_i) = \sum y_i - \bar{y} \sum 1 + b\bar{x} \sum 1 - b \sum x_i = \sum y_i - n\bar{y} + nb\bar{x} - b \sum x_i$$

Now  $\bar{y} = \frac{\sum y_i}{n}$  and  $\bar{x} = \frac{\sum x_i}{n}$  so

$$\sum \varepsilon = \sum y_i - n \frac{\sum y_i}{n} + nb \frac{\sum x_i}{n} - b \sum x_i = \sum y_i - \sum y_i + b \sum x_i - b \sum x_i = 0$$

**b** A linear regression line can be calculated for any set of data, following a linear trend or otherwise, and since by definition the sum of the residuals is equal to zero, this tells you nothing about the actual data itself and so does not guarantee a linear fit. There may be outliers and there is no guarantee that the RSS is close to 0.