# **Correlation and Linear Regression**

$$\sum x^{2} - \frac{(\sum x)^{2}}{n} = 21830 - \frac{360^{2}}{8} = 5630$$

$$\sum y^{2} - \frac{(\sum y)^{2}}{n} = 3500 - \frac{140^{2}}{8} = 1050$$

$$\sum xy - \frac{\sum x \sum y}{n} = 3985 - \frac{(360)(140)}{8} = -2315$$

$$\therefore r = \frac{-2315}{\sqrt{5630 \times 1050}} = -0.952$$

In general, as x increases, y decreases in an almost linear pattern. However, it would be better to sketch a scatter diagram based on the 8 pairs of data to verify.

2(i) 
$$3i$$
 (Using 6C,  $y = -1.121 \times 18.658$  (y on x)  
[LinReg(ax+b).Li, Lx, Y,]  
x values  
Using 6C,  $x = -0.5379 \pm 5.721$  (X on y)x  
(ii)  $2i1$ ) Using line  $d_{1}^{2}$  regression  $d_{2}^{-}$  yon x,  
 $y = -1.121(1-6) \pm 8.658$   
 $= 6.8644$   
Thus there would be  $6864 \approx 6860$  tourists;  
 $2iii$ ) Using line  $d_{2}$  regression  $d_{3} \times 5n g$ ,  
 $x = -0.537(3.2) \pm 5.721$   
 $= 4.00 \times$ 



3(i) 
$$y - \overline{y} = b(x - \overline{x}) \Rightarrow y = bx + (\overline{y} - b\overline{x})$$
  
Since  $y = -0.8x + 13.6$ , by comparing coefficients,  $b = -0.8$ , and  
 $\overline{y} - b\overline{x} = 13.6 \Rightarrow \overline{y} = -0.8(4.5) + 13.6 = 10$   
 $\therefore \Sigma y = 10 \times 8 = 80$   
 $b = -0.8 \Rightarrow \frac{\Sigma xy - \frac{\Sigma x \Sigma y}{n}}{\sum x^2 - \frac{(\Sigma x)^2}{n}} = -0.8$   
 $\therefore \Sigma xy - \frac{\Sigma x \Sigma y}{n} = -0.8 \left[ \Sigma x^2 - \frac{(\Sigma x)^2}{n} \right]$   
 $\Rightarrow \Sigma xy - \frac{36 \times 80}{8} = -0.8 \left[ 204 - \frac{36^2}{8} \right]$   
 $\Rightarrow \Sigma xy = 326.4$   
(ii) Correct:  $\Sigma y = 80$ ,  $\Sigma xy = 326.4$   
From data:  $\Sigma y = 82.7$ ,  $\Sigma xy = 348$   
Difference in  $\Sigma y$  is 21.6  
 $x(2.7) = 21.6 = x = 8$   
9.6 is wrong. When  $x = 8$ ,  $y = 6.9$ .

4(i)	- 0.912
(ii)	Although $r_{\rm B} = -0.912$ is close to $-1$ , which indicate a strong negative linear
	correlation between $X$ and $Y$ , the scatter diagram shows that the data can be better
	represented by a non-linear curve
(iii)	$y = 0.0721(x - 69)^2 + 46.2$
	$r_{(x-69)^2,y} = 0.9981$ ; proposed model is better since  r  is closer to 1 here
(iv)	54.9, As $x = 80$ is out of the range, this estimate is not reliable

5(a)(i)	r = 0.89241 = 0.892 (3 s.f.)
(ii)	r = 0.95956 = 0.960 (3 s.f.)
	Since $r = 0.960$ is closer to 1 than $r = 0.892$ , the model in (i) is less suitable than
	the model in (ii).
(iii)	Regression line of <i>F</i> on <i>x</i> : $F = 0.35903 + 0.029245x$
	F = 0.359 + 0.0292x
	Regression line of x on F: $x = 204.51 + 31.484F$
	x = 205 + 31.5F
(iv)	Using $F = 0.35903 + 0.029245x$ ,
	$100 = 0.35903 + 0.029245t^2$
	t = 58.37047 = 58.4 s (3 s.f.)



(ii)	Based on the scatter diagram, the line 0.906) is quite close to 1. Choose Model (b): $y = ax^b$	ear model is not suitable even though $r$ (=
(11)	Reason: The graph of $y = ax^b$ fits the	scatter diagram better. : From the scatter
	diagram, we see that as $x$ increases,	y increases at a decreasing rate.
(iii)	r = 0.932 $y = ax^{b}$ $\ln y = \ln(ax^{b})$ $\ln y = \ln a + b \ln x$ $\ln y = 4.1912 + 0.3056 \ln x$ $\ln a = 4.1912 \Longrightarrow a = 66.1$ b = 0.306	LinRe9 9=ax+b a=.3055728916 b=4.191193497 r <sup>2</sup> =.8692531131 r=.9323374459
(iv)	when $y = 110$ , $x = 5.29$ . Extrapolation	hence not reliable.

7(i)	Product moment correlation coefficient $r = 0.979$ . There exists a strong positive
	linear correlation between x and y.
(ii)	Equation of least square regression line is $y = 18.5 + 0.564 x$
(iii)	Given that $y = 30$ , $x = 20.4$ from the equation. She left at 7am.
	This estimate is reliable since $r \approx 1$ so we can use the equation of y on x to estimate
	x, giving y and $y = 30$ is within the given data range. (2 reasons)
(iv)	z = time available – time taken
	= 50 - x - y
	=(50-x)-(a+bx)
	= (50 - x) - (18.5 + 0.564x) = 31.5 - 1.564x
(v)	For $z = 0$ , $x = \frac{31.5}{1.564} = 20.14 \approx 20$ min
	The latest time when Ms Chan leaves her house is 7 a.m

8(i)	r = -0.952
	For the depths sampled the moisture content <b>decreases approximately linearly</b>
	with an increase in the depth of the sample.
	Since $r$ is close to 1, this implies that the regression line of $m$ on $x$ and $x$ on $m$ are
	almost identical or close to each other
(ii)	m = -2.2048x + 83.583.
	When $m = 50$ , $x = \frac{83.583 - 50}{2.2048} = 15.232 = 15.2$
	Reliable since the value of <i>m</i> is within the range of the given data. Not
	extrapolating.

9(i)	$PV^c = k$
	$\ln P + c \ln V = \ln k$
	$\ln V = \frac{\ln k}{c} - \frac{1}{c} \ln P$
	$\therefore y = a + bx$ is a straight line, where $a = \frac{\ln k}{c}$ and $b = -\frac{1}{c}$ are constants
(ii)	Using G.C., enter values of P and V in list L1 and L2.
	Let $L3 = \ln P$ and let $L4 = \ln V$
	L1 L2 L3 1 L2 L3 L4 4
	7.7         0         7.7         0         ROLING           2         5.8         .69315         5.8         .69315         1.7579           3         4.5         1.0986         4.5         1.0986         1.5041           4         5         4.5         1.4963         1.2528
	7 2.3 1.9459 2.3 1.9459 .83291 10 1.9 2.3026 1.9 2.3026 .64185 14 1.4 2.6391 1.4 2.6391 .33647
	L100=1 L100=2.041220328

	LinRe9 9=a+bx a=2.143358635 b=6621900222 r <sup>2</sup> =.9877756941
	r = -0.994
(iii)	By using G.C., the required estimated regression line of $y$ on $x$ is
	y = 2.1434 - 0.66219x
	y = 2.14 - 0.662x (to 3 s.f.)
	$\ln V = 2.1434 - 0.66219 \ln P$
	$\ln V = 2.1434 - 0.66219 \ln (8) = 0.76641$
	$V = e^{0.76641} = 2.15203 = 2.2$ (to 1 d.p.) (ans)
	Since $r = -0.994 \approx -1$ which indicates that the sample values of x and y are
	almost perfectly linearly correlated and $\ln(1) \le \ln(8) \le \ln(14)$ , therefore the
	prediction is reliable.
(iv)	$\ln V = \frac{\ln k}{2} - \frac{1}{2} \ln P$
	c = c y = 2.1434 = 0.66219r
	1
	$\frac{1}{c} = 0.66219$
	$c = \frac{1}{10000000000000000000000000000000000$
	0.66219
	$\frac{\ln k}{c} = 2.1434$
	$\ln k = 2.1434 \times 1.51014 = 3.2368$
	$k = e^{3.2368} = 25.452 = 25.5$ (to 3 s.f.)
(v)	For $\sum (y - Y')^2$ to be minimum,
	Then $Y' = a + bx$ must be $y = 2.1434-0.66219x$
	By using G.C., let $L5 = 2.1434 - 0.66219 L3$ and $L6 = (L4 - L5)^2$
	L4     L5     L6     6       2.0412     2.1434     005011     .028156283       1.7579     1.6844     .00577       1.5041     1.4159     .00777       1.2238     1.2254     2.75*6       .83291     .85484     4.85*4       .64185     .61665     5.45*4       .03552     .03552       L6(1)=.0104406851
	Minimum value of $\sum (y - Y')^2 = 0.0282$ (to 3 s.f.)





11(i)	$\overline{x} = 161$ (from calculator or computation)
	when $\bar{x} = 161$ , $\bar{x} = 103.6 + 0.726\bar{y}$
	$\overline{y} = (161 - 103.6) / 0.726$
	= 79.06336088
	using $\overline{y} = \sum y/n$
	$79.06336088 = \frac{1}{6}(65.1 + 73.2 + 85 + k + 80.9 + 89.9)$
	<i>k</i> = 80.3
	Use G.C. to find regression line of <i>y</i> on <i>x</i> :
	y = -97.593 + 1.097x
(ii)	Use <i>y</i> on <i>x</i> line to predict weight.
	When $x = 165$ , $y = -97.593 + 1.097(165)$
	y = 83.4 (1 d.p.) – using 3 d.p. of a and b to compute.
	or
	y = 83.5 - using full accuracy of <i>a</i> and <i>b</i> to compute
(iii)	Using G.C., $r = 0.893$
	y y
	89.9
	85.0
	80.9
	80.3
	73.2
	65.1• x
	150 157 160 162 167 170
	C is unusually overweight.



13(a)(i)	For x on y, $x = -0.3085y + 31.13 \Rightarrow x = -0.309y + 31.1$ (3 s.f.)
	For y on x, $y = -2.8526x + +95.999 \Rightarrow y = -2.85x + 96.0$ (3s.f.)
(ii)	Since chemical Y is the controlled variable, use regression line of x on y.
	$0 = -0.3085y + 31.13 \Longrightarrow y = 100.91$
	The estimation is not valid as this is an extrapolation, linear relation may not hold
	outside the range of data
(b)(i)	By comparing the linear product moment correlation for the 3 models, Model C
	is the most appropriate with the highest value of $ r  = 0.993$ as it best describes
	the data given.

	Using linear transformation $w = \ln x$ , Regression line of w on y is $w = -0.026136y + 3.8294 \Rightarrow w = -0.0261y + 3.83$ (3 s.f.)
(ii)	Change in $w = -0.026136(5) = -0.13068 \approx -0.131 (3 \text{ s.f.})$
	w decreases by 0.131

14(i)	There will be no difference as the product moment correlation coefficient is
	independent of the units in which the data is measured.
(ii)	The <b>regression line of</b> <i>t</i> <b>on</b> <i>x</i> should be used because the running time <i>t</i> is
	<b>dependent</b> on the leg length, <i>x</i>
(iii)	13.90 10.80
	0.70 1.00
(iv)	Yes. Aaron has reason to disagree because the scatter diagram suggests that t and
	x has a <b>curvilinear relationship</b> rather than a linear one.
(v)(a)	Product moment correlation coefficient between t and $\frac{1}{x^2}$ is <b>0.992</b> (3 s.f.)
	The new model is a <b>better model</b> because $ 0.992 $ is <b>closer to 1</b> than
	-0.963  = 0.963.
(b)	Regression line is $t = 7.8603 + 2.8616 \frac{1}{x^2}$ i.e. $t = 7.86 + 2.86 \frac{1}{x^2}$ (3 s.f.)
	when $t = 10$ ,
	$10 = 7.8603 + 2.8616 \frac{1}{x^2}$
	$x^2 = \frac{2.8616}{2.1397}$
	x = 1.16 (to 2 dec places) since $x > 0$
	Thus minimum length of leg required is <b>1.16m</b> .
	This estimate may not be reliable as $t = 10$ is outside the sample data range for
	<i>t</i> .
	OR Extrapolated values are unreliable

15(i)	
(ii)	A linear model is not likely to be appropriate as the area covered would then
	increase continuously, eventually to an infinite area.
(iii)	D = 53: r = -0.99349
	Since $D = 53$ gives a value of r closest to $-1$ , it is the most appropriate
(iv)	Since $D = 55$ gives a value of 7 closest to -1, it is the most appropriate.
$(\mathbf{IV})$	a = 4.20272; b = -0.00899
	Equation of regression line is
	$\ln(D - A) = 4.26272 - 0.60899t$
	When $t = 20$ , $A = 52.99964$ cm <sup>2</sup>
(v)	D is the maximum area of the petri dish

16(i)	A regression line of y on x is more appropriate as the bacteria population depends
10(1)	A regression line of y on x is more appropriate as the bacteria population depends
	on the concentration of nutrients in the water body
(ii)	$r = 0.98119 \approx 0.981 (3 \text{ s.f.})$
(iii)	$4.90 \xrightarrow{y} \\ 19.7 \xrightarrow{x} \\ x \\ 0.101 \\ 0.798 \\ x \\ $
	Although $r \approx 0.981$ suggests a strong linear correlation, the scatter diagram shows that <u>as x increases, y increases at an increasing rate</u> . Therefore, a linear model <u>is</u> <u>not necessarily the best model</u> for the relationship between x and y.
(iv)	By G.C.,
()	$\ln v = 1.3869 + 1.9984 r$
	$1 = 120 \pm 2.00$ (2 - f)
	$\ln y = 1.39 + 2.00x (3 \text{ s.r.})$
	When $x = 1$ ,
	$y = 29.527 \approx 29.5$ (3 s.f.)
	The bacteria population is 29500 (3 s.f.)
	Since $x = 1$ lies outside of the data range, the estimate is not reliable



	$55 = 61.489 - 1.3336\ln(x)$
	x = 129.77 (5sf)
	She will reach a weight of 55kg in the $130^{th}$ month after she started.
	For those who chose $y = c + dx$ in (iii):
	Equation of line as $y = 61.268 - 0.30571x$
	Substituting $y = 55$ into equation to obtain $x = 20.503$
	She will reach a weight of 55kg in the $21^{st}$ month after she started.
(vi)	For any model chosen in (iii):
	As $x \to \infty, y \to -\infty$ .
	Thus, this implies that Amy's weight will decrease to a <b>negative value over time</b> ,
	which is unrealistic.

18(i)	r = -0.9550961661 = -0.955 (to 3 sig.fig)	
	Since r value is close to $-1$ , it suggests a strong negative <u>linear</u> correlation <u>between</u>	
	<u><b>x</b> and </u> <u>y</u> , hence a linear model is appropriate.	
(ii)	v	
20	$\begin{array}{c} & & & \\ & & & & \\ & & & \\ & & & & \\ & & & \\ & & & & & \\ & & & & \\ & & & & & \\ & & & & \\ & & & & \\ & & & & & \\ & & &$	
(iii)	With point P removed, the remaining points lie close to an exponential curve, as $x$	
	increases, y decreases at a decreasing rate, hence consistent with a model of the form	
	$y = Ae^{bx}$ .	
(iv)	$y = Ae^{bx} \Longrightarrow \ln y = \ln A + bx$	
	Since $x$ is the controlled variable, we use the regression line of $\ln y$ on $x$ .	
	From GC, $\ln y = -0.24888x + 3.6276$	
	When $y = 7$ , $\ln 7 = -0.24888x + 3.6276$	
	x = 6.75 = 7 (nearest whole no)	





10(v) 1 day = 24 hours

 $w = -5.31 + 1.35 \ln(24t)$ 

21	(a)(i) Given $y = 79.695x - 205.86$ ,
	$\frac{-}{r}$ 4.07 + 3.84 + 4 + 3.31 + 2.48 + 2.92 + 3.1 + 3.57
	8
	= 3.41125
	$\frac{-}{y} = \frac{140 + 115 + m + 11 + 20 + 25 + 33 + 74}{11 + 20 + 25 + 33 + 74}$
	8
	$=\frac{418+m}{2}$
	Substitute $\overline{x}$ , $\overline{y}$ into the regression line,
	418 + m = (2, 1, 1, 2)
	$\frac{1}{8} = 79.695(3.41125) - 205.86$
	m = 109.997
	m = 110
	(ii) It means that the total fuel use increases by $79.695 \times 10^3$ tonnes when the seat
	capacity increases by 100.
	(b) (i) <sup>y</sup>
	160 (4.07, 140)
	140 -
	120 -
	× ×
	100 -
	80 -
	60 -
	2 2.5 3 3.5 4 4.5
	(ii) The scatter diagram in (b)(i), excluding point $Q$ , suggests that as $x$
	increases, y increases at an increasing rate so model A is not the most
	appropriate.
	(iii) (A): $r = 0.98265$
	(B): $r = 0.9//8/$
	As $ r $ for model (A) is the closest to 1, therefore, model A is the most appropriate
	model.

(iv) Least squares regression line using model A:
$y = 2.6061e^x - 18.429$
When $x = 3.31$ , $y = 2.6061e^{3.31} - 18.429$
= 52.939
So total amount of fuel used is 52 939 tonnes.
(v) Since $r = 0.98265$ is close to 1, which suggests that there is a strong positive
linear relationship between $e^x$ and y, and
x = 3.31 is between 2.48 and 4.07, so interpolation is reliable.

22	Suggested Solution	
(i)		
	t (seconds)	
	115	
	55 $x$ (week number)	
(ii)	A linear model would predict her timing to decrease at a constant rate and eventually negative, which is not possible as there is a limit to how fast a person can swim.	
	A quadratic model would predict that her timings would have a minimum and then increase at an increasing rate, which is also not appropriate.	
(iii)	Based on the scatter diagram and the model, as $x$ increases $t$ decreases at a decreasing rate, therefore $b$ is positive.	
	<i>a</i> has to be positive as it represents the best possible timing that Sharron can swim in the long run.	
(iv)	From GC, r = 0.991 b = 67.69	
	<i>a</i> = 49.50	
(v)	Let <i>m</i> be the best timing Sharron has at the $2^{th}$ month	
(1)	Let m be the best thining sharron has at the of month.	

$\left(\frac{\overline{1}}{x}\right) = 0.33973$	
We know that $\left(\frac{\overline{1}}{x}, \overline{t}\right)$ is on the regression line	
$t = 48.28 + 69.45 \left(\frac{1}{x}\right).$	
$\bar{t} = 48.28 + 69.45(0.33973) = 71.874$	
$\frac{522+m}{8} = 71.874$	
m = 52.992	
Sharron best timing is 53 seconds at the 8th month.	



(b)(iv)
$I = a \mathrm{e}^{bt} \Longrightarrow \ln I = bt + \ln a$
Equation of regression line:
$\ln I = -2.7834239t + 1.6007544 \Longrightarrow \ln I = -2.78t + 1.60$
$\ln a = 1.600754 \Longrightarrow a = 4.96 (3 \text{ s.f.})$
<i>b</i> = -2.78 (3 s.f.)

(b)(v)

t = 0.7, I = 0.706 (to 3 sig fig)

The answer is reliable as *r* is close to -1, and t = 0.7 is within the data range (0.2 to 1.0) and thus the estimate is obtained via interpolation.





(ii) From the scatter diagram (after removing the outlier), as d increases, c decreases at a decreasing rate.

Also, the concentration of the herbicide will not decrease indefinitely and become a negative percentage.

Hence a linear model should not be used to model this set of data.

(iii) Using GC,  $r_A = -0.92958$  while  $r_B = -0.97521$ .

Since the r value for model B is closer to -1 than model A, model B is more appropriate for modelling this set of data.

(iv) 
$$c = ae^{bd}$$
  
 $\ln c = \ln a + bd$   
From GC,  $\ln c = 4.1696 - 0.0066478d$   
 $\ln c = 4.16 - 0.00665d$   
When  $d = 140$ ,  $\ln c = 4.1696 - 0.0066478(140)$ 

 $c=25.5059\approx 25.5$ 

(v) The estimate is unreliable because the data substituted is outside the data range

[20,120] and so the linear relationship between d and  $\ln c$  may not hold.

(vi) Initially, the concentration of herbicides in the soil is 64.7%.