

# **THE ROYAL STATISTICAL SOCIETY**

## **2001 EXAMINATIONS – SOLUTIONS**

### **HIGHER CERTIFICATE**

#### **PAPER I – STATISTICAL THEORY**

The Society provides these solutions to assist candidates preparing for the examinations in future years and for the information of any other persons using the examinations.

The solutions should NOT be seen as "model answers". Rather, they have been written out in considerable detail and are intended as learning aids.

Users of the solutions should always be aware that in many cases there are valid alternative methods. Also, in the many cases where discussion is called for, there may be other valid points that could be made.

While every care has been taken with the preparation of these solutions, the Society will not be responsible for any errors or omissions.

The Society will not enter into any correspondence in respect of these solutions.

Higher Certificate, Paper I, 2001. Question 1

(i)  $A$  and  $B$  are independent. So  $P(\bar{A} \cap \bar{B}) = P(\bar{A})P(\bar{B})$ .

$$P(\bar{A}) = 1 - P(A) = \frac{1}{3}; \quad P(\bar{B}) = \frac{1}{2}; \quad \text{so } P(\bar{A} \cap \bar{B}) = \frac{1}{6}.$$

(ii)  $\bar{A} \cap \bar{B} = [(\bar{A} \cap \bar{B}) \cap C] \cup [(\bar{A} \cap \bar{B}) \cap \bar{C}]$ ,

with the two events  $[(\bar{A} \cap \bar{B}) \cap C]$  and  $[(\bar{A} \cap \bar{B}) \cap \bar{C}]$  being disjoint.

$$\text{Hence } P(\bar{A} \cap \bar{B} \cap \bar{C}) = P(\bar{A} \cap \bar{B}) - P(\bar{A} \cap \bar{B} \cap C) = \frac{1}{6} - \frac{1}{10} = \frac{1}{15}.$$

(iii)  $B \cap C = (A \cap B \cap C) \cup (\bar{A} \cap B \cap C)$ , these being disjoint.

Further,  $(\bar{A} \cap B \cap C) \cup (\bar{A} \cap B \cap \bar{C}) = (\bar{A} \cap B)$ .

$$\begin{aligned} \text{Hence } P(B \cap C) &= P(A \cap B \cap C) + P(\bar{A} \cap B) - P(\bar{A} \cap B \cap \bar{C}) \\ &= \frac{1}{4} + \left(\frac{1}{3} \times \frac{1}{2}\right) - [P(\bar{A} \cap \bar{C}) - P(\bar{A} \cap \bar{B} \cap \bar{C})] \end{aligned}$$

since  $\bar{A}, B$  are independent.

But  $A, C$  are also independent and so are  $\bar{A}, \bar{C}$ .

$$\text{Therefore } P(B \cap C) = \frac{1}{4} + \frac{1}{6} - \left(\frac{1}{3} \times \frac{2}{5}\right) + \frac{1}{15} = \frac{5}{12} - \frac{1}{15} = \frac{63}{12 \times 15} = \frac{7}{20}.$$

$$(iv) \quad P(B|C) = \frac{P(B \cap C)}{P(C)} = \frac{7/20}{3/5} = \frac{7}{12}.$$

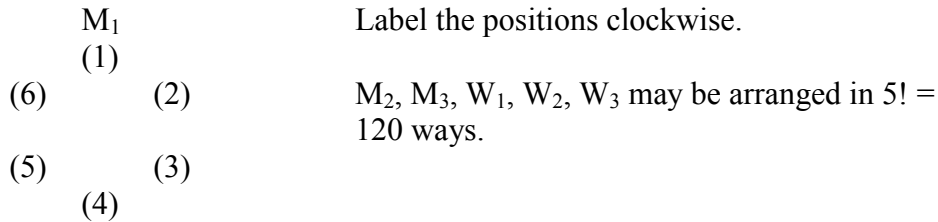
$$(v) \quad P(A|B \cap C) = \frac{P(A \cap B \cap C)}{P(B \cap C)} = \frac{1/4}{7/20} = \frac{5}{7}.$$

$$(vi) \quad P(A \cap B | A \cap C) = \frac{P(A \cap B \cap C)}{P(A \cap C)} = \frac{\frac{1}{4}}{\frac{2}{3} \times \frac{3}{5}} = \frac{5}{8}$$

( $A, C$  are independent).

Higher Certificate, Paper I, 2001. Question 2

- (a) Fix the position of  $M_1$  (suppose him to be the host).



(i)  $M_2, M_3$  must occupy (3) and (5);  $W_1, W_2, W_3$  may occupy the other places in  $3!$  ways, making  $2 \times 3!$  arrangements. The probability is then  $\frac{12}{120} = \frac{1}{10}$ .

(ii)  $M_2, M_3$  must occupy (2) and (6) or (2) and (3) or (5) and (6). In each case,  $M_2, M_3$  can be placed in two orders, making 6 positions altogether for the three men. The women may again fill the remaining places in  $3!$  ways.

The probability is  $\frac{6 \times 6}{120} = \frac{3}{10}$ .

(iii) EITHER  $1 - \frac{1}{10} - \frac{3}{10} = \frac{3}{5}$ , because this is the only other arrangement possible besides (i) and (ii);

OR by having  $M_2$  in (2),  $M_3$  in (4) or (5);  $M_2$  in (6),  $M_3$  in (3) or (4);  $M_2$  in (3),  $M_3$  in (4);  $M_2$  in (4),  $M_3$  in (5); or any of these with  $M_2, M_3$  interchanged, giving 12 positionings of the men. There are again  $3!$  orders for the women, so the probability is  $\frac{6 \times 12}{120} = \frac{3}{5}$ .

- (b) (i) Event  $D$  is "has disease",  $T$  is "tests positive".

$$\begin{aligned}
 P(D|T) &= \frac{P(D \cap T)}{P(T)} = \frac{P(T|D)P(D)}{P(T)} = \frac{P(T|D)P(D)}{P(T|D)P(D) + P(T|\bar{D})P(\bar{D})} \\
 &= \frac{p_1 p_0}{p_1 p_0 + (1 - p_2)(1 - p_0)}.
 \end{aligned}$$

(ii)  $\frac{0.95 \times 0.005}{(0.95 \times 0.005) + (0.05 \times 0.995)} = \frac{0.00475}{0.0545} = 0.0872$ .

The error rates in the clinical tests are large compared to the chance of having the disease, so the calculated probability is very small.

Higher Certificate, Paper I, 2001. Question 3

(a) (i) 
$$P(S \geq 2300) = 1 - P(S < 2300) = 1 - \Phi\left(\frac{2300 - 2000}{300}\right)$$
$$= 1 - \Phi(1) = 1 - 0.8413 = 0.1587.$$

$$P(H \geq 2300) = 1 - \Phi\left(\frac{2300 - 2500}{125}\right) = 1 - \Phi(-1.6)$$
$$= 1 - 0.0548 = 0.9452.$$

(ii)  $P(S > H) = P(S - H > 0)$ , where  $(S - H)$  is  $N(-500, 90000 + 15625)$   
i.e.  $N(-500, 325^2)$ . Hence  $P(S > H) = \Phi\left(\frac{-500}{325}\right) = \Phi(-1.5385) = 0.0620$ .

(b) (i) The lifetime  $X$  is  $S$  with probability 0.6 and  $H$  with probability 0.4.

Hence  $E[X] = 0.6 \times 2000 + 0.4 \times 2500 = 2200$  hrs.

$$P(X > 2600) = P(X > 2600 | S)P(S) + P(X > 2600 | H)P(H)$$
$$= \Phi\left(-\frac{600}{300}\right) \times 0.6 + \Phi\left(-\frac{100}{125}\right) \times 0.4$$

(using the appropriate tail areas from Normal tables)

$$= 0.6 \Phi(-2) + 0.4 \Phi(-0.8) = 0.6 \times 0.02275 + 0.4 \times 0.2119$$
$$= 0.01365 + 0.08476 = 0.09841.$$

(ii) 
$$P(H | > 2600) = \frac{P(X > 2600 | H)P(H)}{P(X > 2600)} = \frac{0.08476}{0.09841} = 0.8613.$$

(iii)  $\bar{X}$  will have mean 2200. It is not Normally distributed but we may apply the Central Limit Theorem if we know its variance. In large samples we may take  $\bar{X}$  as approximately  $N(2200, 346.8^2)$ , so that

$$P(\bar{X} > 2300) = 1 - \Phi\left(\frac{2300 - 2200}{346.8/\sqrt{100}}\right) = 1 - \Phi\left(\frac{100}{34.68}\right) = \Phi(-2.8835) = 0.002.$$

Higher Certificate, Paper I, 2001. Question 4

An easy method is to consider  $X$  as  $\sum X_i$ , where  $X_i$  are a set of  $n$  Bernoulli variables with  $P(X_i = 1) = p$ ,  $P(X_i = 0) = (1 - p)$ . Then  $E[X_i] = p$ , so  $E[X] = np$ .

Also  $E[X_i^2] = p$ , so  $\text{Var}(X_i) = p - p^2$  and  $\text{Var}(X) = n(p - p^2) = npq$ .

ALTERNATIVELY: 
$$E[X] = \sum_{x=0}^n x \binom{n}{x} p^x q^{n-x} = \sum_{x=1}^n \frac{n! p^x q^{n-x}}{(x-1)!(n-x)!}$$
$$= np \sum_{x=1}^n \binom{n-1}{x-1} p^{x-1} q^{(n-1)-(x-1)} = np.$$

Similarly,  $\text{Var}(X) = E[X(X-1)] + E[X] - (E[X])^2$ , and we have

$$E[X(X-1)] = \sum_{x=0}^n x(x-1) \binom{n}{x} p^x q^{n-x} = \sum_{x=2}^n x(x-1) \binom{n}{x} p^x q^{n-x}$$
$$= n(n-1) p^2 \sum_{x=2}^n \binom{n-2}{x-2} p^{x-2} q^{(n-2)-(x-2)} = n(n-1) p^2,$$

and hence  $\text{Var}(X) = n(n-1) p^2 + np - n^2 p^2 = np - np^2 = npq$ .

PGFs or MGFs could also be used.

(i) (a)  $0.75^{16} \approx 0.0100226 = 0.0100$  approx.

(b)  $1 - P(\text{no one gets all 16 right})$ , probability is  $1 - \{1 - 0.75^{16}\}^{12}$   
 $= 1 - \{0.9899774\}^{12} = 0.1139$ .

(ii)  $P(B - A > 0)$  can be studied using a Normal approximation to the difference  $B - A$ , i.e.  $N(16\{0.5 - 0.75\}, 16\{(0.5 \times 0.5) + (0.75 \times 0.25)\})$ , i.e.  $N(-4, 7)$ .

The probability is found as  $P\left(B - A > \frac{1}{2}\right)$  using a continuity correction since  $B - A$  takes discrete values.

Hence it is  $1 - \Phi\left(\frac{0.5 - (-4)}{\sqrt{7}}\right) = \Phi\left(-\frac{4.5}{\sqrt{7}}\right) = \Phi(-1.7008) \approx 0.0445$ .

[Note: this would be 0.0653 without the continuity correction.]

(iii)  $E[\bar{X}] = E[X] = np = 16 \times 0.75 = 12$  in set  $A$ .

Similarly,  $E[\bar{Y}] = 16 \times 0.5 = 8$  in set  $B$ .

There are 12 students in  $A$  and 25 in  $B$ , so that

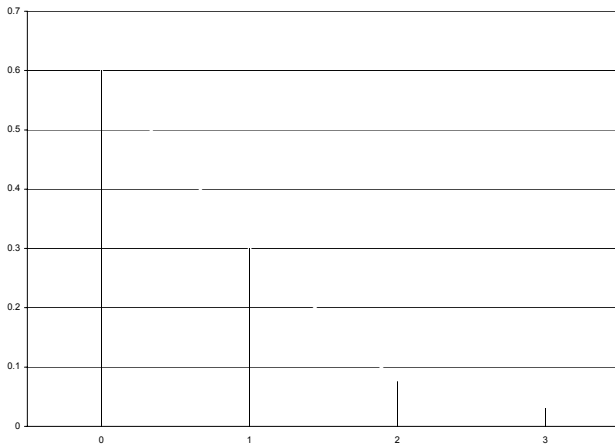
$$\text{Var}(\bar{X}) = \frac{16 \times 0.75 \times 0.25}{12} = \frac{1}{4} \quad \text{in set } A$$

$$\text{Var}(\bar{Y}) = \frac{16 \times 0.5 \times 0.5}{25} = \frac{4}{25} \quad \text{in set } B.$$

Higher Certificate, Paper I, 2001. Question 5

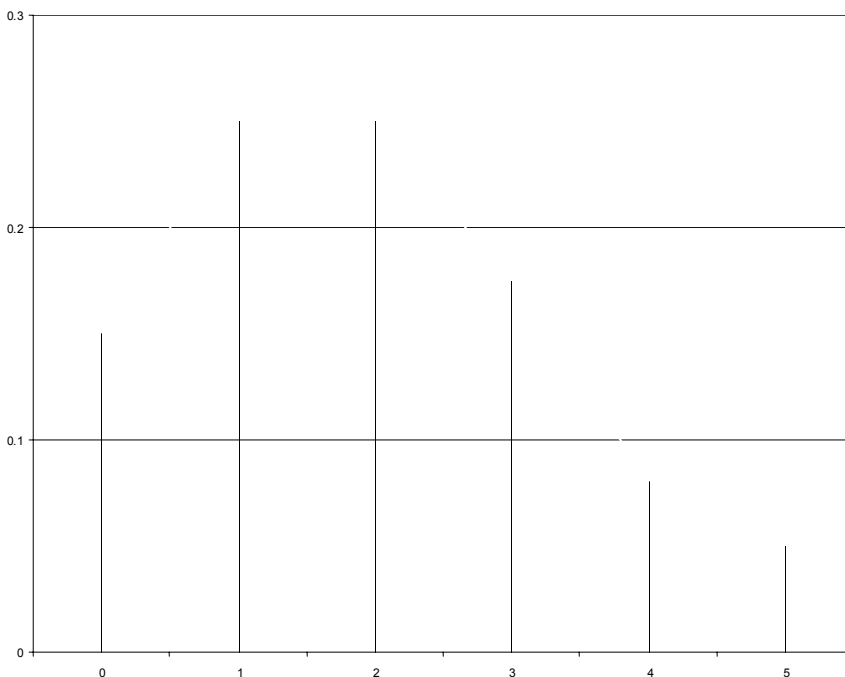
$$\frac{p(x+1)}{p(x)} = \frac{e^{-\lambda} \lambda^{x+1}}{(x+1)! e^{-\lambda} \lambda^x} = \frac{\lambda}{x+1} \text{ for } x = 0, 1, 2, \dots$$

For  $\lambda = \frac{1}{2}$ ,  $p(0) = 0.60653$ ; so  $p(1) = 0.30327$ ,  $p(2) = 0.07582$ ,  $p(3) = 0.01264$ .



Graph of  $p(x)$  for  $\lambda = \frac{1}{2}$ .

For  $\lambda = 2$ ,  $p(0) = 0.13534$ ;  
so  $p(1) = 0.27067 = p(2)$ ,  $p(3) = 0.18045$ ,  $p(4) = 0.09022$ ,  $p(5) = 0.03609$ .



Graph of  $p(x)$  for  $\lambda = 2$ .

$$M_X(t) = E[e^{Xt}] = \sum_{x=0}^{\infty} \frac{e^{xt} e^{-\lambda} \lambda^x}{x!} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{(\lambda e^t)^x}{x!} = e^{-\lambda} \cdot e^{\lambda e^t} = \exp(\lambda\{e^t - 1\}).$$

$$\frac{\partial M}{\partial t} = \lambda e^t e^{\lambda(e^t-1)}; \text{ put } t=0 \text{ and this is } \lambda, \text{ which is therefore } E[X].$$

$$\frac{\partial^2 M}{\partial t^2} = (\lambda^2 e^{2t} + \lambda e^t) e^{\lambda(e^t-1)}; \text{ put } t=0 \text{ and this is } \lambda^2 + \lambda, \text{ but it is also } E[X^2].$$

$$\text{Hence } \text{Var}(X) = E[X^2] - (E[X])^2 = \lambda^2 + \lambda - (\lambda)^2 = \lambda.$$

$$\frac{\partial^3 M}{\partial t^3} = (\lambda^3 e^{3t} + 3\lambda^2 e^{2t} + \lambda e^t) e^{\lambda(e^t-1)} = \lambda^3 + 3\lambda^2 + \lambda \text{ at } t=0. \text{ This is } E[X^3].$$

$$\begin{aligned} \text{Now, } E[(X-\lambda)^3] &= E[X^3] - 3\lambda E[X^2] + 3\lambda^2 E[X] - \lambda^3 \\ &= \lambda^3 + 3\lambda^2 + \lambda - 3\lambda(\lambda^2 + \lambda) + 3\lambda^2 \cdot \lambda - \lambda^3 = \lambda. \end{aligned}$$

For  $Y = X_1 + X_2 + \dots + X_n$  we have

$$E[e^{Yt}] = \prod_{i=1}^n E[e^{X_i t}] = (M_X(t))^n = \exp[\lambda n(e^t - 1)].$$

This is the mgf of a Poisson distribution with parameter  $\lambda n$  and so  $E[Y] = \text{Var}(Y) = \lambda n$ .

$$P(Y \geq 40) = 1 - P(Y \leq 39) \approx 1 - \Phi\left(\frac{39.5 - 25}{5}\right),$$

using continuity correction, and  $\mu = \lambda n = 25$ . This is  $1 - \Phi(2.9) = 0.00187$ .

With a positively skew distribution, the Normal approximation is likely to underestimate the probability in the right hand tail and so we expect this answer to be less than the true value.



Higher Certificate, Paper I, 2001. Question 6

$$L(p) = \prod_{i=1}^n (q^{x_i-1} p) = p^n q^{\sum x_i - n} = \left(\frac{p}{1-p}\right)^n (1-p)^{\bar{n}}$$

$$\ln L = n \ln p - n \ln(1-p) + n\bar{x} \ln(1-p)$$

$$\frac{\partial}{\partial p}(\ln L) = \frac{n}{p} + \frac{n}{1-p} - \frac{n\bar{x}}{1-p} = 0 \quad \text{when} \quad \frac{1}{\hat{p}} = \frac{-1+\bar{x}}{1-\hat{p}} \quad \text{which gives} \quad \hat{p} = \frac{1}{\bar{x}} \quad \text{as m.l. estimate.}$$

$$\frac{\partial^2(\ln L)}{\partial p^2} = -\frac{n}{p^2} - \frac{n(\bar{x}-1)}{(1-p)^2} \quad \text{which is} < 0, \quad \text{confirming the maximum.}$$

$$E\left[\frac{\partial^2 \ln L}{\partial p^2}\right] = \frac{n}{p^2} + \frac{n}{(1-p)^2} [E(\bar{X}) - 1] = \frac{n}{p^2} + \frac{n}{(1-p)^2} \left(\frac{1}{p} - 1\right) = \frac{n}{p^2} + \frac{n}{p(1-p)}$$

$$= \frac{n}{p^2(1-p)}$$

$$\text{Hence } \text{Var}(\hat{p}) \approx \frac{p^2(1-p)}{n}.$$

$$\sum fx = 448, \quad \sum f = 56, \quad \hat{p} = \frac{56}{448} = 0.125.$$

$$\text{Var}(\hat{p}) = \frac{0.125^2 \times 0.875}{56} = 0.0002441, \quad \text{SE}(\hat{p}) = 0.015625.$$

Approximate 95% confidence interval for  $p$  is  $\hat{p} \pm 1.96 \text{SE}(\hat{p})$ , which is  
 $0.125 \pm 1.96 \times 0.015625$ , i.e.  $0.125 \pm 0.030625$  or  $(0.0944, 0.1556)$ .

When  $p = \frac{1}{6}$ , we have  $\frac{1}{\bar{X}} \sim N\left(\frac{1}{6}, \frac{5}{216 \times 56}\right)$ , i.e.  $N(0.1667, 0.00041336)$ ; therefore

the probability of obtaining  $\hat{p} \leq 0.125$  is approximately

$$\Phi\left(\frac{0.125 - 0.1667}{0.020331}\right) = \Phi(-2.0496) = -0.0202.$$

The confidence interval for  $p$  did not include  $\frac{1}{6}$ ; also now the probability being very small is consistent with rejecting a null hypothesis that  $p = \frac{1}{6}$ , i.e. that the die is fair.

Higher Certificate, Paper I, 2001. Question 7

(a)  $P(X_i \leq x) = F(x)$  for  $i = 1, 2, \dots, n$ .

(i) 
$$F_{X_{\max}}(x) = P(X_1 \leq x, X_2 \leq x, \dots, X_n \leq x) = \prod_{i=1}^n P(X_i \leq x)$$
$$= [F(x)]^n$$

(ii) For  $X_{\min} \geq x$ , we require every  $X_i$  to be  $\geq x$ .

Now,  $P(X_i \leq x) = F(x)$ , so  $P(X_i \geq x) = 1 - F(x)$ , for all  $i$ .

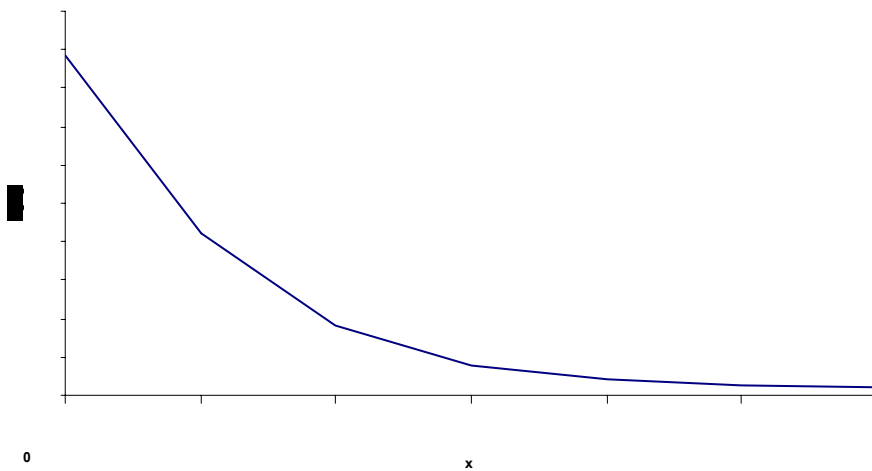
$$\begin{aligned} \therefore F_{X_{\min}}(x) &= P(X_{\min} \leq x) = 1 - P(X_{\min} \geq x) = 1 - P(X_1 \geq x, \dots, X_n \geq x) \\ &= 1 - [1 - F(x)]^n \end{aligned}$$

(iii) Pdfs are derivatives of cdfs:

$$f_{X_{\max}} = n[F(x)]^{n-1} f(x)$$

$$f_{X_{\min}} = n[1 - F(x)]^{n-1} f(x), \quad \text{since } \frac{\partial F(x)}{\partial x} = f(x).$$

(b)



$$F(x) = \int_0^x \frac{\alpha}{(1+u)^{\alpha+1}} du = \left[ -\frac{1}{(1+u)^\alpha} \right]_0^x = 1 - \frac{1}{(1+x)^\alpha}.$$

Median  $M$  is such that  $F(M) = \frac{1}{2}$ . So we have

$$1 - \frac{1}{(1+M)^\alpha} = \frac{1}{2} \quad \text{or} \quad \frac{1}{(1+M)^\alpha} = \frac{1}{2} \quad \text{or} \quad 2 = (1+M)^\alpha \quad \text{or} \quad M = 2^{(1/\alpha)} - 1.$$

Using (a)(ii),  $F_{X_{\min}} = 1 - \left( \frac{1}{(1+x)^\alpha} \right)^n = 1 - \frac{1}{(1+x)^{n\alpha}}$ , also Pareto but with  $\alpha$  replaced

by  $n\alpha$ .

The median of  $X_{\min}$  is then  $2^{\frac{1}{n\alpha}} - 1$ , which is  $2^{1/n} - 1$  if  $\alpha = 1$ . We require

$$2^{1/n} - 1 < 0.1 \quad \text{or} \quad 2^{1/n} < 1.1, \quad \text{i.e.} \quad \frac{1}{n} \ln 2 < \ln 1.1, \quad \text{giving} \quad n > \frac{\ln 2}{\ln 1.1} = \frac{0.6931}{0.0953} = 7.27.$$

Hence  $n \geq 8$ .

Higher Certificate, Paper I, 2001. Question 8

$$y_i = \alpha + \beta x_i + \varepsilon_i$$

where  $y$  is the response (observation) and  $x$  the value of the explanatory variable on that unit;  $\{\varepsilon_i\}$  is a set of independent, identically distributed random variables with mean 0 and the same variance  $\sigma^2$ . Usually they are assumed Normal as a basis for inference.  $x_i$  is assumed "fixed", not "random".

(i) (a) There is an increasing trend, and the relationship between  $y$  and  $x$  appears curvilinear.

(b) In simple regression  $R^2$  is the square of the correlation,  $r$ , between  $x$  and  $y$ . In general it is the proportion of the variance of  $y$  which can be explained by the dependence of  $y$  on all explanatory variables  $\{x_i\}$  in the model; hence it is the square of the correlation between  $\hat{y}$  and  $y$ .

In the ANOVA,  $\frac{\text{regression SS}}{\text{total SS}} = \frac{269.33}{335.37} = 0.803$ , or 80.3%.

(c) As  $x$  (% operating capacity) increases by 1 unit so  $y$  (profit) increases by 0.31562 units.

A 95% confidence interval is  $0.31562 \pm t_{10} \times 0.04942$ , which is  $0.31562 \pm 0.11011$  or (0.2055, 0.4257).

(d) Values of profit have been predicted for capacity 25%, 50% and 75%. The confidence intervals for these predictions are those given; but note that 25% is far outside the range of available data (hence the remark about extreme  $x$  values). A 95% confidence interval is an interval which should cover the true  $y$  at a given  $x$  with probability 0.95, based on the fitted linear regression.

(ii) (a) The logarithmic plot shows that a linear regression in these units is a much better fit. There is still an increasing trend.

(b)  $\log_{10}(\text{profits}) = -0.519 + 0.0177(\text{capacity})$

i.e.  $\text{profits} = 10^{-0.519+0.0177(\text{capacity})}$ .

(c) For capacity = 25, the 95% limits are  $-0.2731$  and  $+0.1210$ , and the actual prediction is  $-0.0760$ , in  $\log_{10}$  units. Anti-logging these (i.e. raising 10 to these powers) we find the 95% limits are 0.5332 and 1.3213.

The 'prediction' is 0.8395.

These limits do not overlap the limits on the previous model.

The prediction now is for a small profit, compared with a loss on the previous model.

- (iii) The scatter plots indicate that the logarithmic model is preferred, and so do the plots of residuals which show a random pattern (as compared with a systematic, curved one for the previous model).  $R^2$  also higher (92.4%) on the log model.

But a log model cannot predict negative profits – i.e. losses – which are quite possible in general though not for these data if used within the range of  $x$  values given.

Extrapolation down to 25% is well outside the data and so is not reliable on any model.