

THE ROYAL STATISTICAL SOCIETY

2004 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, 2004. Question 1

- (i) (a) $P(\text{all four favour the complex}) = (0.6)^4$.
 $P(\text{all four oppose the complex}) = (0.3)^4$.
 $P(\text{all four are indifferent}) = (0.1)^4$.
- So $P(\text{all four think alike}) = (0.6)^4 + (0.3)^4 + (0.1)^4 = 0.1378$.
- (b) $P(\text{an individual is not opposed}) = 0.6 + 0.1 = 0.7$.
- So $P(\text{none of the four is opposed}) = (0.7)^4 = 0.2401$.
- (c) Possible favourable results are *FFOI*, *FOOI*, *FOII*, in any order.
- $P(\text{FFOI}) = (0.6)^2(0.3)(0.1) = 0.0108$
 $P(\text{FOOI}) = (0.6)(0.3)^2(0.1) = 0.0054$
 $P(\text{FOII}) = (0.6)(0.3)(0.1)^2 = 0.0018$
- Each result can be arranged in $\frac{4!}{2!1!1!} = 12$ ways.
- So overall probability is $12(0.0108 + 0.0054 + 0.0018) = 0.216$.
- (d) From (a), $P(\text{all four in favour}) = (0.6)^4 = 0.1296$. From (b), $P(\text{none opposed}) = 0.2401$. So the required conditional probability is
- $0.1296/0.2401 = 0.5398$.
- (ii) The number in favour is binomially distributed with $n = 4$ and $p = 0.6$. So the expectation (mean) is $4 \times 0.6 = 2.4$ and the variance is $4 \times 0.6 \times 0.4 = 0.96$.
- (iii) $P(\text{opposed}) = P(\text{opposed} | \text{young})P(\text{young}) + P(\text{opposed} | \text{older})P(\text{older})$
 $= (0.12 \times 0.25) + (p \times 0.75)$
- where $p = P(\text{opposed} | \text{older})$. But we are given that $P(\text{opposed}) = 0.3$. Hence $p = 0.36$.
- (iv) In samples of one "young" and three "olders",
- $P(\text{exactly one opposes}) = P(\text{"young" opposes, "olders" do not})$
 $+ P(\text{"young" does not oppose, one "older" opposes})$
 $= \{(0.12)(0.64)^3\} + \{3(0.88)(0.36)(0.64)^2\} = 0.03146 + 0.38928 = 0.4207$.

Higher Certificate, Paper I, 2004. Question 2

- (i) (a) The moment generating function is

$$M_X(t) = E(e^{tX}) = \sum_{x=0}^{\infty} e^{tx} \frac{e^{-\lambda} \lambda^x}{x!} = \sum_{x=0}^{\infty} e^{-\lambda} \frac{(\lambda e^t)^x}{x!} = \exp(\lambda e^t - \lambda) = \exp(\lambda(e^t - 1))$$

(b) $E(X) = M_X'(0) = \left[\lambda e^t e^{\lambda(e^t - 1)} \right]_{t=0} = \lambda.$

$$E(X^2) = M_X''(0) = \left[\frac{d}{dt} \left(\lambda e^t e^{\lambda(e^t - 1)} \right) \right]_{t=0}$$
$$= \left[\lambda e^t \cdot e^{\lambda(e^t - 1)} \lambda e^t + e^{\lambda(e^t - 1)} \cdot \lambda e^t \right]_{t=0} = \lambda^2 + \lambda.$$

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

(Alternatively, the moments could be obtained from the power series expansion of $M_X(t)$.)

(Alternatively, though with comparatively lengthy algebra, the moments could be obtained directly by $E(X) = \sum xP(X=x)$ and $E(X^2) = \sum x^2P(X=x)$; or, somewhat easier, use $E[X(X-1)] = \sum x(x-1)P(X=x)$ (this is λ^2) and then $\text{Var}(X) = E[X(X-1)] + E(X) - \{E(X)\}^2$.)

- (c) The binomial distribution with parameters n and p may be approximated by the Poisson distribution with parameter np if n is large and p is small. As a "rule of thumb", $\frac{1}{2} \leq np \leq 10$ gives an indication of how large n should be and how small p should be. (If $np > 10$, a Normal approximation to the binomial may be better.)

- (ii) Let X = number of wrong calculations. We have $X \sim B(200, 0.0075)$.

$$P(X=1) = \binom{200}{1} (0.0075)(0.9925)^{199} = 200 \times 0.0075 \times 0.9925^{199} = 0.3353(2).$$

$$P(X=4) = \binom{200}{4} (0.0075)^4 (0.9925)^{196}$$
$$= \frac{200 \times 199 \times 198 \times 197}{4 \times 3 \times 2 \times 1} \times 0.0075^4 \times 0.9925^{196} = 0.0468(0).$$

Continued on next page

(iii) We approximate using $X \sim \text{Poisson}(200 \times 0.0075 = 1.5)$. With this,

$$P(X = 1) = 1.5e^{-1.5} = 0.3347(0)$$

giving a percentage error of $\frac{100(0.33532 - 0.33470)}{0.33532} = 0.18\%$, and

$$P(X = 4) = \frac{e^{-1.5} (1.5)^4}{4!} = 0.0470(7)$$

giving a percentage error of $\frac{100(0.04707 - 0.04680)}{0.04680} = 0.58\%$.

[Note. These percentage errors might come out *slightly* differently if more accuracy is kept in the binomial and Poisson probabilities.]

Both approximations are remarkably accurate, with percentage errors well below 1%. The approximation for $X = 1$ (one wrong calculation) is the more accurate of the two. That approximation is an underestimate; the other is an overestimate.

Higher Certificate, Paper I, 2004. Question 3

Actual volume $X \sim N(1010, 8^2)$. Let $Z \sim N(0,1)$.

(i)
$$P(X < 1000) = P\left(Z < \frac{1000 - 1010}{8}\right) = P(Z < -1.25) = 0.1056.$$

(ii) Let Y be the total volume in a 6-pack.

We have $Y \sim N(6 \times 1010, 64 + 64 + 64 + 64 + 64 + 64)$, i.e. $Y \sim N(6060, 384)$.

$$P(Y < 6000) = P\left(Z < \frac{6000 - 6060}{\sqrt{384}}\right) = P(Z < -3.06) = 0.0011.$$

(Alternatively, could use $\bar{X} \sim N(1010, 64/6)$ and calculate $P(\bar{X} < 1000)$.)

This probability is considerably smaller than that in part (i). In practical terms, this is because there will be a tendency for heavier and lighter cartons in a 6-pack to balance each other out. Alternatively, in terms of probability distributions, consider X and \bar{X} : \bar{X} has the same mean as X but only one-sixth of the variance, so less of the lower tail of the distribution of \bar{X} is below the nominal volume of 1000.

(iii) The new volume $W \sim N(\mu, 4^2)$, where μ is the new mean. So we have that $P(W < 1000) = P\left(Z < \frac{1000 - \mu}{4}\right)$. We require that this probability must be no greater than 0.1056. Thus the cut-off point for Z is to be $z = -1.25$ (as before). Hence $\frac{1000 - \mu}{4} = -1.25$, giving $\mu = 1005$.

This means that 5 ml per carton could be saved, i.e. a cost saving per carton of $\frac{5}{1000} \times \text{£}1$. To recover the £200, the number of cartons required is $\frac{200}{5/1000} = 40000$.

Higher Certificate, Paper I, 2004. Question 4

(i) The binomial distribution with parameters n, p can be approximated by $N(np, np(1-p))$ when n is large and p is not too near to 0 or 1. As a "rule of thumb", the approximation is likely to be good if both np and $np(1-p)$ are > 10 .

Let $X \sim B(n, p)$ and let Φ denote the c.d.f. of $N(0, 1)$. Using a continuity correction,

$$P(X \leq x) \approx \Phi\left(\frac{x + \frac{1}{2} - np}{\sqrt{np(1-p)}}\right) \text{ and } P(X < x) \approx \Phi\left(\frac{x - \frac{1}{2} - np}{\sqrt{np(1-p)}}\right).$$

The 95% confidence interval for p uses the estimated variance $\hat{p}(1-\hat{p})/n$, giving the approximate interval

$$\hat{p} - 1.96\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} < p < \hat{p} + 1.96\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

The estimate of p is $\hat{p} = \frac{30}{50} = 0.6$, so $\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \sqrt{\frac{(0.6)(0.4)}{50}} = 0.0693$. Thus the approximate interval is

$$0.6 - (1.96 \times 0.0693), \quad 0.6 + (1.96 \times 0.0693)$$

i.e. (0.464, 0.736).

(ii) $P(\text{neither hits}) = (1-p)^2$. Therefore $P(\text{at least 1 hit}) = 1 - (1-p)^2 = p(2-p)$. We estimate this by $(0.6)(2-0.6) = 0.84$.

(iii) When $p = 0.464$ (lower limit of interval in part (i)), we have $p(2-p) = 0.713$. Similarly, when $p = 0.736$, we have $p(2-p) = 0.930$. Thus (0.713, 0.930) is the required interval.

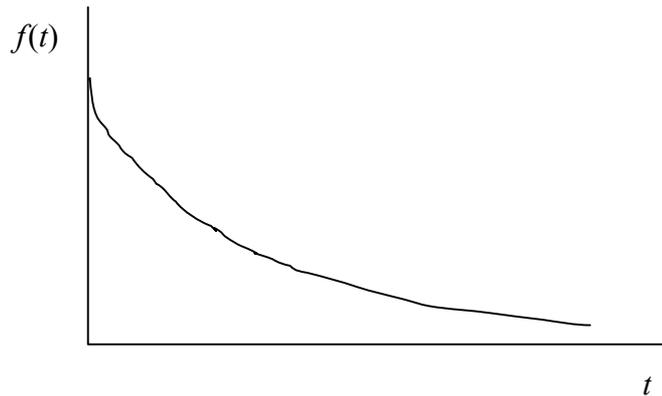
(iv) When n pairs are fired, $P(\text{all miss}) = [(1-p)^2]^n$, estimated by $(0.16)^n$. Hence $(0.16)^n < 0.0005$ is required. Solving this by taking logarithms to base 10, we have $n \log_{10}(0.16) < \log_{10}(0.0005)$, i.e. $-0.79588n < -3.30103$ which gives $n > 4.148$. so n must be at least 5.

Higher Certificate, Paper I, 2004. Question 5

(i) $f(t) = \lambda e^{-\lambda t}, \quad t > 0; \lambda > 0$

(a) Sketch of $f(t)$.

[**NOTE.** The curve should of course appear as a smooth decaying exponential; it might not do so, due to the limits of electronic reproduction.]



(b) C.d.f. is $F(t) = P(T \leq t) = \int_0^t \lambda e^{-\lambda v} dv = \lambda \left[-\frac{1}{\lambda} e^{-\lambda v} \right]_0^t = 1 - e^{-\lambda t}$.

(c) $P(a < T \leq b) = F(b) - F(a) = e^{-\lambda a} - e^{-\lambda b}$.

(ii) Assume all settlements of invoices are independent.

$P(50 \text{ in first week}) = \{F(1)\}^{50} = (1 - e^{-\lambda})^{50}$, because $T \leq 1$ for all these 50.

Likewise, $1 < T \leq 2$ for the 35 in the second week, so we have $P(35 \text{ in second week}) = \{F(2) - F(1)\}^{35} = (e^{-\lambda} - e^{-2\lambda})^{35}$.

The remaining 15 have $T > 2$, which has probability $1 - P(T \leq 2) = e^{-2\lambda}$, and thus $P(15 \text{ after week 2}) = (e^{-2\lambda})^{15}$.

The likelihood is therefore the product

$$L(\lambda) = k(1 - e^{-\lambda})^{50} (e^{-\lambda} - e^{-2\lambda})^{35} (e^{-2\lambda})^{15}$$

where k is a constant of proportionality.

Continued on next page

Taking logarithms (base e),

$$\begin{aligned}\log L(\lambda) &= \log k + 50 \log(1 - e^{-\lambda}) + 35 \log\{e^{-\lambda}(1 - e^{-\lambda})\} + 15 \log(e^{-2\lambda}) \\ &= \log k + 85 \log(1 - e^{-\lambda}) - (35 + 30)\lambda = \log k + 85 \log(1 - e^{-\lambda}) - 65\lambda.\end{aligned}$$

$$\therefore \frac{d}{d\lambda} \log L = \frac{85e^{-\lambda}}{1 - e^{-\lambda}} - 65 = \frac{85}{e^{\lambda} - 1} - 65.$$

Equating to zero, $85 = 65(e^{\lambda} - 1)$ or $e^{\lambda} = 150/65$, so that $\hat{\lambda} = \log(150/65) = 0.836$.

[It is easy to check that this is indeed a maximum; e.g. $\frac{d^2}{d\lambda^2} \log L = -\frac{85}{(e^{\lambda} - 1)^2} < 0$.]

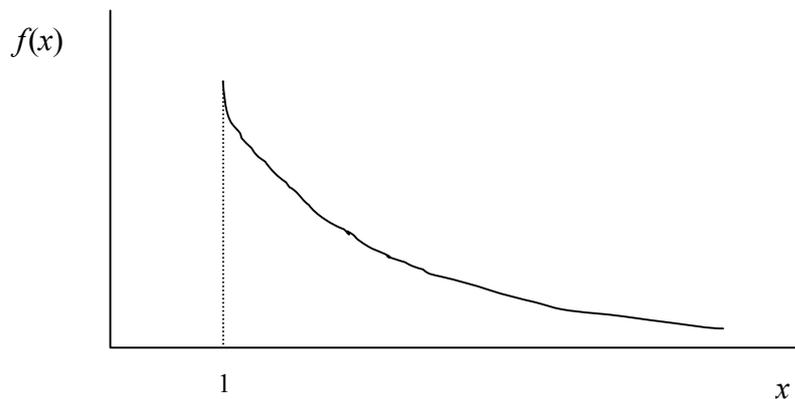
(iii) $1 - e^{-0.836} = 0.5666$; $e^{-0.836} - e^{-1.672} = 0.43344 - 0.18787 = 0.2456$. Hence, out of 100 invoices, 56.66, 24.56 and 18.78 would be expected to be paid, on this model, in weeks 1, 2 and later. The actual numbers were 50, 35 and 15. The prediction for the second week is a long way from what happened, balanced by smaller discrepancies in the other two periods. This does not seem very satisfactory.

Higher Certificate, Paper I, 2004. Question 6

$$f(x) = \frac{k}{x^{k+1}}, \quad x \geq 1; \quad k > 0$$

(i) Sketch of $f(x)$.

[NOTE. The curve should of course appear as a smooth curve; it might not do so, due to the limits of electronic reproduction.]



$$\text{C.d.f. is } F(x) = P(X \leq x) = \int_1^x \frac{k}{u^{k+1}} du = \left[-\frac{1}{u^k} \right]_1^x = 1 - \frac{1}{x^k} \quad (\text{for } x \geq 1).$$

(ii) Median M has $\frac{1}{2} = F(M) = 1 - M^{-k}$, so $\frac{1}{2} = M^{-k}$ and hence $M = 2^{1/k}$.

Lower quartile Q_1 has $\frac{1}{4} = F(Q_1) = 1 - Q_1^{-k}$, so $\frac{3}{4} = Q_1^{-k}$, i.e. $Q_1 = (4/3)^{1/k}$.

Upper quartile Q_3 has $\frac{3}{4} = F(Q_3) = 1 - Q_3^{-k}$, so $Q_3 = (4)^{1/k}$.

Hence the semi-interquartile range is $\frac{1}{2} \left\{ 4^{1/k} - \left(\frac{4}{3} \right)^{1/k} \right\}$.

Continued on next page

$$(iii) \quad E(X) = \int_1^{\infty} xf(x) dx = \int_1^{\infty} \frac{k}{x^k} dx = \left[\frac{-k}{(k-1)x^{k-1}} \right]_1^{\infty} = \frac{k}{k-1}.$$

$$E(X^2) = \int_1^{\infty} x^2 f(x) dx = \int_1^{\infty} \frac{k}{x^{k-1}} dx = \left[\frac{-k}{(k-2)x^{k-2}} \right]_1^{\infty} = \frac{k}{k-2}.$$

$$\therefore \text{Var}(X) = E(X^2) - \{E(X)\}^2 = \frac{k}{k-2} - \frac{k^2}{(k-1)^2}$$

$$= \frac{k}{(k-2)(k-1)^2} \{(k-1)^2 - k(k-2)\} = \frac{k}{(k-1)^2(k-2)}.$$

$$P(X > E(X)) = \int_{k/(k-1)}^{\infty} \frac{k}{x^{k+1}} dx = \left[-\frac{1}{x^k} \right]_{k/(k-1)}^{\infty} = \left(\frac{k-1}{k} \right)^k, \text{ or this can be written down directly from the c.d.f. found in part (i).}$$

(iv) For the case $k = 3$,

(a) $M = 2^{1/3}$ in the units given, or £12599,

(b) mean = $3/2$ in the units given, or £15000,

(c) inserting $X = 10$, $P(X \leq 10) = 1 - \frac{1}{10^3}$, so $P(X > 10) = \frac{1}{10^3}$, i.e. 0.1%.

Higher Certificate, Paper I, 2004. Question 7

$$(i) \quad E(X) = \int_0^\theta \frac{x}{\theta} dx = \frac{1}{\theta} \left[\frac{1}{2} x^2 \right]_0^\theta = \frac{1}{2} \theta.$$

$$E(X^2) = \int_0^\theta \frac{x^2}{\theta} dx = \frac{1}{\theta} \left[\frac{1}{3} x^3 \right]_0^\theta = \frac{1}{3} \theta^2.$$

$$\therefore \text{Var}(X) = E(X^2) - \{E(X)\}^2 = \frac{1}{3} \theta^2 - \left(\frac{1}{2} \theta \right)^2 = \frac{1}{12} \theta^2.$$

$$(ii) \quad P(\text{longest offcut is } \leq x) = P(\text{all } n \text{ offcuts are } \leq x).$$

The c.d.f. for each X_i is $F(x) = P(X \leq x) = \int_0^x \frac{du}{\theta} = \left[\frac{u}{\theta} \right]_0^x = \frac{x}{\theta}$, and the X_i are all independent. Therefore $P(\text{all } n \text{ offcuts are } \leq x) = \{F(x)\}^n = \left(\frac{x}{\theta} \right)^n$, and this is also $P(\text{longest offcut is } \leq x)$, i.e. the c.d.f. of the sample maximum $X_{(n)}$. Thus the p.d.f. of $X_{(n)}$ is the derivative of this, i.e. nx^{n-1}/θ^n . This is for the interval $(0, \theta)$.

$$\therefore E(X_{(n)}) = \int_0^\theta \frac{nx^n}{\theta^n} dx = \frac{n}{\theta^n} \left[\frac{x^{n+1}}{n+1} \right]_0^\theta = \frac{n\theta}{n+1}.$$

$$E(X_{(n)}^2) = \int_0^\theta \frac{nx^{n+1}}{\theta^n} dx = \frac{n}{\theta^n} \left[\frac{x^{n+2}}{n+2} \right]_0^\theta = \frac{n\theta^2}{n+2}.$$

$$\begin{aligned} \therefore \text{Var}(X_{(n)}) &= E(X_{(n)}^2) - \{E(X_{(n)})\}^2 = \frac{n\theta^2}{n+2} - \frac{n^2\theta^2}{(n+1)^2} \\ &= n\theta^2 \left(\frac{(n+1)^2 - n(n+2)}{(n+2)(n+1)^2} \right) = \frac{n\theta^2}{(n+1)^2(n+2)}. \end{aligned}$$

Immediately we have $E\left(\frac{n+1}{n} X_{(n)}\right) = \theta$, so $\frac{n+1}{n} X_{(n)}$ is an unbiased estimator of θ .

$$\text{Var}\left(\frac{n+1}{n} X_{(n)}\right) = \frac{(n+1)^2}{n^2} \text{Var}(X_{(n)}) = \frac{(n+1)^2}{n^2} \frac{n\theta^2}{(n+1)^2(n+2)} = \frac{\theta^2}{n(n+2)}.$$

Continued on next page

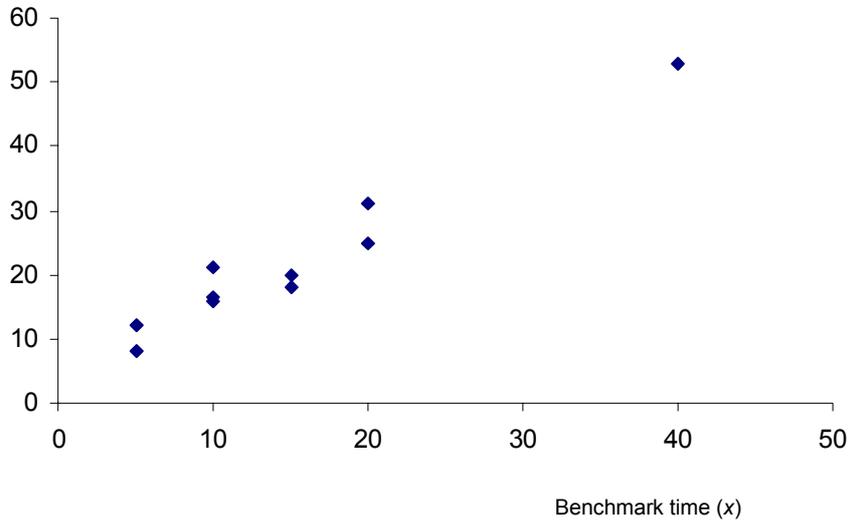
(iii) We have (see part (i)) that $E(X) = \theta/2$. Thus the method of moments estimator of $\theta/2$ is \bar{X} , and so the method of moments estimator of θ is $2\bar{X}$ or $\frac{2}{n}\sum X_i$ as required.

$$\text{Var}\left(\frac{2}{n}\sum X_i\right) = \text{Var}(2\bar{X}) = 4\text{Var}(\bar{X}) = \frac{4}{n}\text{Var}(X) = \frac{4}{n} \cdot \frac{\theta^2}{12} = \frac{\theta^2}{3n}.$$

Higher Certificate, Paper I, 2004. Question 8

(i)

Trainee's time (y)



Simple linear regression analysis seems quite suitable.

(ii) The model is $y_i = \alpha + \beta x_i + e_i$, where $\{e_i\}$ are uncorrelated with zero mean and (constant) variance σ^2 (independent identically distributed $N(0, \sigma^2)$ for the purpose of undertaking statistical tests, as in part (iii)). Estimating by the method of least squares gives

$$\hat{\beta} = \frac{S_{xy}}{S_{xx}}, \quad \hat{\alpha} = \bar{y} - \hat{\beta}\bar{x},$$

where (standard notation)

$$S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i y_i - \frac{\sum x_i \sum y_i}{n},$$

$$S_{xx} = \sum (x_i - \bar{x})^2 = \sum x_i^2 - \frac{(\sum x_i)^2}{n}.$$

We have

$$\hat{\beta} = \frac{S_{xy}}{S_{xx}} = \frac{4440 - (150 \times 220 / 10)}{3200 - (150^2 / 10)} = \frac{1140}{950} = 1.20 \quad \text{and} \quad \hat{\alpha} = 22 - (1.20 \times 15) = 4,$$

so the line is

$$y = 4 + 1.2x.$$

Continued on next page

The total sum of squares is $S_{yy} = \sum (y_i - \bar{y})^2 = \sum y_i^2 - \frac{(\sum y_i)^2}{10} = 1440$.

The sum of squares for regression is $\hat{\beta}S_{xy}$ (or S_{xy}^2 / S_{xx}) = 1368.

Therefore the residual sum of squares is $1440 - 1368 = 72$.

This has 8 degrees of freedom, so the residual mean square ($\hat{\sigma}^2$) is $72/8 = 9$.

The coefficient of determination $R^2 = 1368/1440 = 0.95$ (usually given as 95%).

(iii) The estimated variance of $\hat{\beta}$ is $9/950 = 0.009474$. So the test statistic for testing the null hypothesis $\beta = 1$ is $\frac{1.2 - 1}{\sqrt{0.009474}} = 2.05$, which we refer to t_8 .

This is not significant at the 5% level, so the null hypothesis $\beta = 1$ cannot be rejected.

(iv) The model here is $y_i = bx_i + e_i$.

Estimating b by least squares, we minimise $\Omega = \sum_{i=1}^n (y_i - bx_i)^2$.

Differentiating with respect to b , we have $\frac{d\Omega}{db} = -2 \sum (y_i - bx_i) x_i$.

Setting this equal to zero gives $\sum x_i y_i = \hat{b} \sum x_i^2$, i.e. $\hat{b} = \sum x_i y_i / \sum x_i^2$.

(Note that $\frac{d^2\Omega}{db^2} = 2 \sum x_i^2 > 0$, so this is a minimum.)

Thus we have $\hat{b} = 4440/3200 = 1.3875$.