

**THE ROYAL STATISTICAL SOCIETY
2015 EXAMINATIONS – SOLUTIONS
GRADUATE DIPLOMA – MODULE 5**

The Society is providing these solutions to assist candidates preparing for the examinations in 2017.

The solutions are intended as learning aids and should not be seen as "model answers".

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.

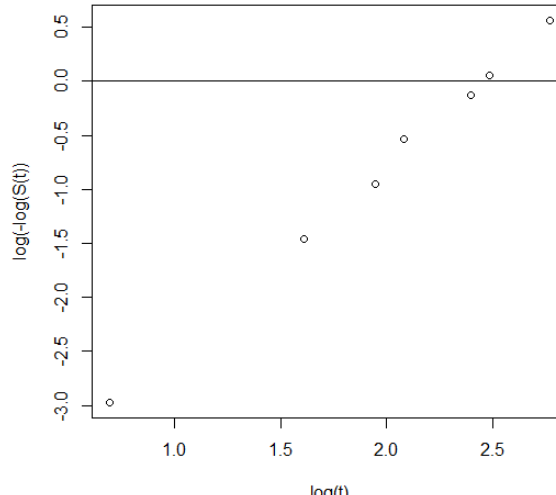
Question 1

(i)	The plots suggest that subjects with higher Sodium and Potassium levels and lower Carbon Dioxide levels are more likely to be Type B, but there is not much difference in Age between the groups [1 mark]. Some plots show good separation between the groups, notably Potassium vs Age and Potassium vs Carbon Dioxide [1 mark]. It appears that (some of) these variables will be useful for classification [1 mark].
(ii)	With leaving-one-out cross-validation, each subject in turn is omitted from the dataset, and the discriminant function estimated from the remaining subjects' data is used to predict the class of the omitted one [1 mark]. The misclassification rates are estimated by comparing the true class of each subject with the predicted class when it is left out [1 mark]. This procedure reduces the downward bias in misclassification results that is found when the discriminant is estimated and used on the same dataset [1 mark]. In this example, linear discrimination seems likely to be very useful since the misclassification rates obtained, even with cross-validation, are low [1 mark]. Type B individuals are particularly well classified [1 mark].
(iii)	LD-A = 2254.04 [1 mark]. LD-B = 2254.93 [1 mark]. Since LD-A < LD-B, classify the new individual as Type B [1 mark].
(iv)	Linear discriminant analysis (LDA) assumes that the classification variables jointly have a multivariate normal distribution [1 mark], separately within the populations of Type A and Type B individuals [1 mark], though the procedure is robust against minor departures from this assumption [1 mark]. The population covariance matrices are assumed to be equal [1 mark]. Quadratic discriminant analysis (QDA) does not require the population covariance matrices to be equal [1 mark].
(v)	Without cross-validation, the results with QDA seem better than with LDA [1 mark]. With cross-validation, however, the results with QDA are worse than with LDA especially for Type B individuals [1 mark]. Recommend using LDA [1 mark], since the cross-validated misclassification rates give a better guide to how the discriminants will perform with new patients [1 mark].

Question 2

(i)	The purpose of cluster analysis is to determine whether a set of objects may be meaningfully divided into different classes [1 mark] and, if so, to determine those classes [1 mark] using similarities and dissimilarities between the objects [1 mark].
(ii)	<p>Euclidean distance: $\left(\sum_{i=1}^p (x_i - y_i)^2\right)^{1/2}$ [1 mark].</p> <p><i>Candidates should define and discuss two of the following, for 2 marks each.</i></p> <p>Squared Euclidean distance = $\sum_{i=1}^p (x_i - y_i)^2$ [1 mark]. This will have the effect of increasing the distances between pairs of observations, making them less likely to cluster together and hence increase the number of clusters [1 mark].</p> <p>Manhattan distance = $\sum_{i=1}^p x_i - y_i$ [1 mark]. This will have the effect of decreasing the distances between pairs of observations, making them more likely to cluster together and hence decrease the number of clusters [1 mark].</p> <p>Chebyshev distance = $\max_i x_i - y_i$ [1 mark]. This will have the effect of increasing the distances between pairs of observations, making them less likely to cluster together and hence increase the number of clusters [1 mark].</p> <p>Mahalanobis distance = $\left([\underline{\mathbf{x}} - \underline{\mathbf{y}}]^T \mathbf{S}^{-1} [\underline{\mathbf{x}} - \underline{\mathbf{y}}]\right)^{1/2}$ [1 mark]. It is not clear in general what effect this will have on the number of clusters as that will depend on the pattern of correlations between variables [1 mark].</p>
(iii)	Both methods of clustering begin with all objects in their own cluster [1 mark]. At the first step, pairs of objects whose distance is less than a certain threshold value are then combined into a new cluster [1 mark] and this process continues until no distances between clusters are below the threshold [1 mark]. In single-linkage, at later steps the distance between two clusters is determined by the distance of the two closest objects (nearest neighbours) in the different clusters [1 mark]. In complete-linkage, the distance between two clusters is determined by the greatest distance between any two objects in the different clusters [1 mark]. As a result, single-linkage will usually give more clusters than complete-linkage [1 mark].
(iv)	k means clustering starts by randomly generating k clusters [1 mark], determining the location of the cluster centre then assigning each point to the cluster whose centre is closest to it [1 mark], iteratively re-computing the new cluster centres until convergence occurs [1 mark], which is identified by point-cluster assignments no longer changing [1 mark]. The idea is to move objects between clusters with the goal of minimising variability within clusters and maximising variability between clusters. A major disadvantage of this method is that k must be set before the cluster analysis is started [1 mark], but it is often computationally faster than hierarchical methods [1 mark].

Question 3

(i)	$S(t) = \int_t^{\infty} \lambda \gamma t^{\gamma-1} e^{-\lambda t^\gamma} dt \quad [1 \text{ mark}]$ $= \int_{\lambda t^\gamma}^{\infty} e^{-u} du \quad \text{setting } u = \lambda t^\gamma \quad [1 \text{ mark}]$ $= \left[-e^{-u} \right]_{\lambda t^\gamma}^{\infty} = e^{-\lambda t^\gamma} \quad [1 \text{ mark}]$																																																
(ii)	$\log S(t) = \log(e^{-\lambda t^\gamma}) = -\lambda t^\gamma \quad [1 \text{ mark}]$ $\log[-\log S(t)] = \log[\lambda t^\gamma] \quad [1 \text{ mark}]$ $= \log(\lambda) + \gamma \log(t) \quad [1 \text{ mark}]$ <p>So, the intercept would be $\log(\lambda)$ [1 mark] and the slope γ [1 mark].</p>																																																
(iii)	<table style="width: 100%; border-collapse: collapse; text-align: center;"> <thead> <tr> <th style="padding: 5px;">t_i</th> <th style="padding: 5px;">n_i</th> <th style="padding: 5px;">d_i</th> <th style="padding: 5px;">$\frac{n_i - d_i}{n_i}$</th> <th style="padding: 5px;">$\hat{S}(t_i)$</th> <th style="padding: 5px;">$\log\{-\log[\hat{S}(t_i)]\}$</th> </tr> </thead> <tbody> <tr><td style="padding: 5px;">2</td><td style="padding: 5px;">20</td><td style="padding: 5px;">1</td><td style="padding: 5px;">19/20</td><td style="padding: 5px;">0.950</td><td style="padding: 5px;">-2.97</td></tr> <tr><td style="padding: 5px;">5</td><td style="padding: 5px;">18</td><td style="padding: 5px;">3</td><td style="padding: 5px;">15/18</td><td style="padding: 5px;">0.792</td><td style="padding: 5px;">-1.46</td></tr> <tr><td style="padding: 5px;">7</td><td style="padding: 5px;">14</td><td style="padding: 5px;">2</td><td style="padding: 5px;">12/14</td><td style="padding: 5px;">0.679</td><td style="padding: 5px;">-0.949</td></tr> <tr><td style="padding: 5px;">8</td><td style="padding: 5px;">11</td><td style="padding: 5px;">2</td><td style="padding: 5px;">9/11</td><td style="padding: 5px;">0.555</td><td style="padding: 5px;">-0.530</td></tr> <tr><td style="padding: 5px;">11</td><td style="padding: 5px;">8</td><td style="padding: 5px;">2</td><td style="padding: 5px;">6/ 8</td><td style="padding: 5px;">0.416</td><td style="padding: 5px;">-0.131</td></tr> <tr><td style="padding: 5px;">12</td><td style="padding: 5px;">6</td><td style="padding: 5px;">1</td><td style="padding: 5px;">5/ 6</td><td style="padding: 5px;">0.347</td><td style="padding: 5px;">0.057</td></tr> <tr><td style="padding: 5px;">16</td><td style="padding: 5px;">2</td><td style="padding: 5px;">1</td><td style="padding: 5px;">1/ 2</td><td style="padding: 5px;">0.173</td><td style="padding: 5px;">0.562</td></tr> </tbody> </table> <p>[The final column of the table above is not required until part (iv). ½ mark for method + ½ mark for each correct row leading to a correct value of $S(t)$.]</p>	t_i	n_i	d_i	$\frac{n_i - d_i}{n_i}$	$\hat{S}(t_i)$	$\log\{-\log[\hat{S}(t_i)]\}$	2	20	1	19/20	0.950	-2.97	5	18	3	15/18	0.792	-1.46	7	14	2	12/14	0.679	-0.949	8	11	2	9/11	0.555	-0.530	11	8	2	6/ 8	0.416	-0.131	12	6	1	5/ 6	0.347	0.057	16	2	1	1/ 2	0.173	0.562
t_i	n_i	d_i	$\frac{n_i - d_i}{n_i}$	$\hat{S}(t_i)$	$\log\{-\log[\hat{S}(t_i)]\}$																																												
2	20	1	19/20	0.950	-2.97																																												
5	18	3	15/18	0.792	-1.46																																												
7	14	2	12/14	0.679	-0.949																																												
8	11	2	9/11	0.555	-0.530																																												
11	8	2	6/ 8	0.416	-0.131																																												
12	6	1	5/ 6	0.347	0.057																																												
16	2	1	1/ 2	0.173	0.562																																												
(iv)	<div style="text-align: center;">  </div>																																																

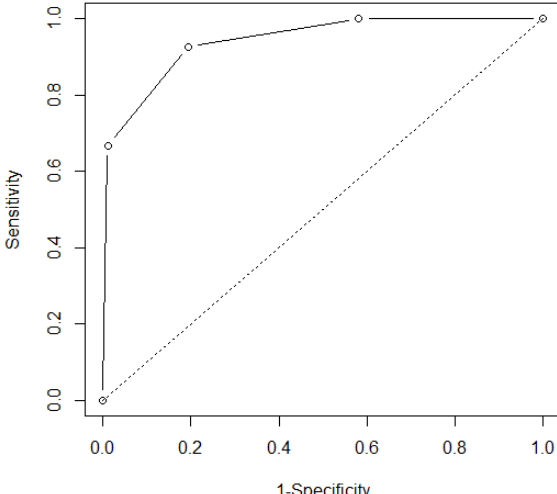
Question 3

(iv)	<p>... continued</p> <p>[1 mark for correct method of calculating $\log(-\log(S(t)))$ values, 1 mark for correct values, 1 mark for appropriate graph, 1 mark for correct points on graph]</p> <p>The points lie approximately on a straight line, so a Weibull model seems appropriate [1 mark].</p>
(v)	<p>An estimate of γ, the slope, would be about $3.5/2.5 = 1.4$ [1 mark]. An estimate of the intercept, $\log(\lambda)$, would be about $-3.0 - 0.7 = -3.7$ [1 mark], giving an estimate of λ of about 0.025 [1 mark].</p>

Question 4

(i)	<p>The model is:</p> $h(t; \underline{z}) = h_0(t) \cdot \exp\{\underline{\beta}^T \underline{z}\} = h_0(t) \cdot \exp\{\beta_1 z_1 + \beta_2 z_2 + \beta_3 z_3 + \beta_4 z_4\} \quad [1 \text{ mark}]$ <p>Here, $h_0(t)$ is the baseline hazard function for subjects who are given standard care, are male and not white, and have 0% burns [1 mark]. z_1 is an indicator variable for Treatment, so β_1 summarises the average effect of being treated with additional care rather than standard care. z_2 is an indicator variable for Gender, so β_2 summarises the average effect of being Female rather than Male. z_3 is an indicator variable for Race, so β_3 summarises the effect of being White rather than Non White. <i>Idea of an indicator variable, 1 mark; idea of an average difference in effect, 1 mark.</i> z_4 is the percentage burns so β_4 summarises the average effect of each additional % of body surface area affected [1 mark]. <i>Alternative, equivalent parameterisations should also be given full marks.</i></p>
(ii)	<p>$b_1/\text{se}(b_1) = -0.606/0.296 = -2.05 < -1.96$, so Treatment is a useful explanatory variable [1 mark]. Since the sign of b_1 is negative, additional care will reduce the hazard of infection (all other factors being equal) [1 mark].</p> <p>$b_2/\text{se}(b_2) = -0.631/0.390 = -1.62 > -1.96$, so Gender is not a useful explanatory variable [1 mark]. There is no significant difference on average between the hazard for Male and Female subjects (all other factors being equal) [1 mark].</p> <p>$b_3/\text{se}(b_3) = 2.12/1.01 = 2.10 > 1.96$, so Race is a useful explanatory variable [1 mark]. Since the sign of b_3 is positive, White subjects have a higher hazard of infection (all other factors being equal) [1 mark].</p> <p>$b_4/\text{se}(b_4) = 0.00404/0.00703 = 0.57 < 1.96$, so Severity is not a useful explanatory variable [1 mark]. There is no significant difference on average between the hazard for subjects with different initial severity levels (all other factors being equal) [1 mark].</p>
(iii)	<p>The hazard ratio is $\exp\{\beta_1\}$ [1 mark]. A 95% CI for β_1 is:</p> $b_1 \pm 1.96 \text{ se}(b_1), \text{ i.e. } -0.606 \pm (1.96 \times 0.296), \text{ i.e. } (-1.186, -0.026) \quad [1 \text{ mark}]$ <p>Therefore, a 95% CI for the hazard ratio is (0.306, 0.974) [1 mark].</p>
(iv)	<p>Add 3 new indicator variables to the model, corresponding to any 3 of the causes [1 mark]. For example, one suitable variable would be coded 1 if the cause was chemical and 0 otherwise [1 mark].</p>
(v)	<p>This might have been done because there were so few instances of each of the other causes, or because there were clinical reasons for thinking that the effects of the other kinds of burn would be similar [1 mark]. The effect of the new variable is not significant: $-0.481/0.330 = -1.46 > -1.96$ [1 mark].</p>

Question 5

(i)	<table style="margin-left: auto; margin-right: auto; border-collapse: collapse;"> <thead> <tr> <th style="border: none;"></th> <th style="border: none; text-align: center; border-bottom: 1px solid black;">True +</th> <th style="border: none; text-align: center; border-bottom: 1px solid black;">True -</th> </tr> </thead> <tbody> <tr> <td style="border: none; padding-right: 10px;">Test +</td> <td style="border: none; text-align: center;">18</td> <td style="border: none; text-align: center;">1</td> </tr> <tr> <td style="border: none; padding-right: 10px;">Test -</td> <td style="border: none; text-align: center;">9</td> <td style="border: none; text-align: center;">92</td> </tr> </tbody> </table> <p style="margin-left: 40px;">Sensitivity = $\frac{18}{27} = 0.667$ [1 mark] Specificity = $\frac{92}{93} = 0.989$ [1 mark] PPV = $\frac{18}{19} = 0.947$ [1 mark] NPV = $\frac{92}{101} = 0.911$ [1 mark]</p> <p>PPV and NPV are both high, meaning that patients classified as positive or negative for hypothyroidism are very likely to be positive or negative respectively [1 mark], which is important for patients [1 mark]. Specificity is high, meaning that almost all patients who are truly negative are classified as negative by this test [1 mark]. Sensitivity is only moderate so almost one-third of truly positive patients will be classified as negative and this makes the test less useful than it might be [1 mark].</p>		True +	True -	Test +	18	1	Test -	9	92
	True +	True -								
Test +	18	1								
Test -	9	92								
(ii)	<p>$c = 7$: sensitivity = $\frac{25}{27} = 0.926$, specificity = $\frac{75}{93} = 0.806$ [2 marks] $c = 9$: sensitivity = $\frac{27}{27} = 1$, specificity = $\frac{39}{93} = 0.419$ [2 marks]</p> <div style="text-align: center;">  </div> <p>[1 mark for correct method, 1 mark for correct axis labels, 1 mark for correct points, 1 mark for $y=x$ line]</p> <p>Recommend $c \approx 7$ [1 mark] as that achieves the best balance between sensitivity and specificity [1 mark] in the absence of information about costs of misdiagnosis and prevalence.</p>									
(iii)	<p>The “ideal” ROC curve has AUC = 1 [1 mark], so a value of 0.86 suggests that this procedure is potentially useful for diagnosis [1 mark].</p>									

Question 6

(i)	<p>In a case-control study, cases (observed to have an outcome of interest) are identified, along with appropriate controls (observed to be free of the outcome) [1 mark]. Enquiries are then made to find out whether or not each subject in each group had previously been exposed to a certain risk factor, and the frequency of exposure is compared between the groups [1 mark]. One advantage of case-control studies relative to prospective studies (such as cohort studies) is that results can be obtained more quickly, without having to wait for some subjects to develop the outcome [1 mark]. Case-control studies are particularly beneficial for studying rare outcomes, since other types of studies would need a large cohort to be enrolled in order to generate sufficient outcomes for analysis [1 mark].</p>												
(ii)	<p>The main advantage of a matched case-control study is an improvement in the precision for estimating effect sizes (or an increase in power for statistical tests) [1 mark]. This arises because the subjects in a matched pair are similar in respect of key variables and the matching reduces residual variance [1 mark]. A disadvantage is that it can be difficult to find controls who exactly match each case [1 mark] and this can increase the cost or extend the time required to carry out the study [1 mark].</p>												
(iii)	<p>For the unmatched analysis, require the following table.</p> <table style="margin-left: auto; margin-right: auto; border-collapse: collapse;"> <tr> <td style="padding: 5px;"></td> <td style="text-align: center; padding: 5px;"><u>Smoker</u></td> <td style="text-align: center; padding: 5px;"><u>Non-Smoker</u></td> <td style="padding: 5px;"></td> </tr> <tr> <td style="text-align: right; padding: 5px;">Control</td> <td style="text-align: center; padding: 5px;">37</td> <td style="text-align: center; padding: 5px;">130</td> <td style="padding: 5px;"></td> </tr> <tr> <td style="text-align: right; padding: 5px;">Case</td> <td style="text-align: center; padding: 5px;">55</td> <td style="text-align: center; padding: 5px;">112</td> <td style="padding: 5px;">[1 mark for table, possibly implied]</td> </tr> </table> <p style="margin-left: 20px;">Unmatched OR = $\frac{55/112}{37/130} = 1.73$ [1 mark]</p> <p style="margin-left: 20px;">Matched OR = $\frac{40}{22} = 1.82$ [1 mark for method, 1 mark for correct value]</p>		<u>Smoker</u>	<u>Non-Smoker</u>		Control	37	130		Case	55	112	[1 mark for table, possibly implied]
	<u>Smoker</u>	<u>Non-Smoker</u>											
Control	37	130											
Case	55	112	[1 mark for table, possibly implied]										
(iv)	<p>H_0: there is no association between mother's smoking behaviour and low birth weight [1 mark]</p> <p>For the unmatched analysis, expected frequencies are: 46, 121, 46, 121 [1 mark].</p> <p>So: $\chi^2 = 9^2 \left(\frac{2}{46} + \frac{2}{121} \right) = 4.86$ [1 mark for formula, 1 mark for correct final value]</p> <p>Since $\chi^2 > 3.841$ reject H_0 at 5% level [1 mark].</p> <p>For McNemar's Test, $\chi^2 = \left(\frac{ 40-22 -1}{40+22} \right)^2 = 4.66$ [1 mark for formula, 1 mark for correct final value]</p> <p>Since $\chi^2 > 3.841$ reject H_0 at 5% level [1 mark].</p>												

Question 7

(i)	<p>Consider any element, x', in the population. Number of different possible samples that include $x' = \binom{N-1}{n-1}$ [1 mark]. Each of these samples has probability $\binom{N}{n}$ of being chosen. So, probability that x' is chosen = $\frac{\binom{N-1}{n-1}}{\binom{N}{n}} = \frac{n}{N}$ [1 mark].</p>
(ii)	$E(\bar{x}) = E\left\{\frac{1}{n} \sum_{i=1}^n x_i\right\} = \sum_{j=1}^N \frac{x_j}{n} P(x_j \text{ in sample}) \quad [1 \text{ mark}]$ $= \sum_{j=1}^N \frac{x_j}{n} \frac{n}{N} = \frac{1}{N} \sum_{j=1}^N x_j = \mu \quad [1 \text{ mark}]$
(iii)	$\text{Var}(\bar{x}) = \frac{1}{n^2} \text{Var}\left\{\sum_{i=1}^n x_i\right\} \quad [1 \text{ mark}]$ $= \frac{1}{n^2} \{n\sigma^2 + n(n-1)\tau\}, \text{ where } \tau = \text{Cov}(x', x) \quad [1 \text{ mark}]$ <p>Now $\tau = E(x'x) - \mu^2 \quad [1 \text{ mark}]$</p> <p>and $E\{(x')^2\} = \sigma^2 + \mu^2 \quad [1 \text{ mark}]$</p> <p>But $E\{(x')^2\} = E\left\{x' \left[N\mu - \sum_{x_j \neq x'} x_j \right]\right\}$</p> $= N\mu^2 - (N-1)E(x'x) \quad [1 \text{ mark}]$ <p>So $\sigma^2 + \mu^2 = N\mu^2 - (N-1)E(x'x)$</p> <p>i.e. $E(x'x) = \mu^2 - \sigma^2 / (N-1) \quad [1 \text{ mark}]$</p> <p>i.e. $\tau = -\sigma^2 / (N-1) \quad [1 \text{ mark}]$</p> $\text{Var}(\bar{x}) = \frac{1}{n^2} \left\{ n\sigma^2 - \frac{n(n-1)\sigma^2}{N-1} \right\} = \frac{\sigma^2}{n} \frac{N-n}{N-1} \quad [1 \text{ mark}]$

Question 7

<p>(iv)</p>	<p>$2\sqrt{\text{Var}(\bar{x})} = B$</p> <p>i.e. $2\sqrt{\frac{\sigma^2}{n} \frac{N-n}{N-1}} = B$ [1 mark]</p> <p>i.e. $\frac{\sigma^2}{n} \frac{N-n}{N-1} = \frac{B^2}{4}$ [1 mark]</p> <p>i.e. $\sigma^2 N - \sigma^2 n = \frac{B^2}{4} (N-1)n$</p> <p>i.e. $\sigma^2 N = \left\{ \frac{B^2}{4} (N-1) + \sigma^2 \right\} n$ [1 mark]</p> <p>i.e. $n = \frac{4\sigma^2 N}{B^2 (N-1) + 4\sigma^2}$ [1 mark]</p>
<p>(v)</p>	<p>This is an example of part (iv) with $N = 1000$, $\sigma^2 = 100$, $B = 2$ [1 mark].</p> <p>So: $n = \frac{4 \times 100 \times 1000}{(4 \times 999) + (4 \times 100)} = 90.99$ [1 mark substitution, 1 mark value]</p> <p>Need to sample at least 91 chickens. [1 mark]</p>

Question 8

(i)	<table style="margin-left: auto; margin-right: auto; border-collapse: collapse;"> <thead> <tr> <th style="border-bottom: 1px solid black;">i</th> <th style="border-bottom: 1px solid black;">N_i</th> <th style="border-bottom: 1px solid black;">σ_i</th> <th style="border-bottom: 1px solid black;">c_i</th> <th style="border-bottom: 1px solid black;">$N_i\sigma_i/\sqrt{c_i}$</th> </tr> </thead> <tbody> <tr><td>1</td><td>1000</td><td>1550</td><td>1.5</td><td>1265570</td></tr> <tr><td>2</td><td>1100</td><td>2150</td><td>1.5</td><td>1931014</td></tr> <tr><td>3</td><td>1600</td><td>2950</td><td>1.5</td><td>3853864</td></tr> <tr><td>4</td><td>600</td><td>1100</td><td>1</td><td>660000</td></tr> <tr><td>Total</td><td>4300</td><td></td><td></td><td>7710448</td></tr> </tbody> </table> <p style="margin-top: 10px;">For proportional allocation, $n_i \approx \frac{N_i}{N}n$ [1 mark], giving 116, 128, 186, 70 [1 mark].</p> <p>For optimal allocation, $n_i \approx \frac{N_i\sigma_i/\sqrt{c_i}}{\sum N_i\sigma_i/\sqrt{c_i}}n$ [1 mark], giving 82, 125, 250, 43 [1 mark for $N_i\sigma_i/\sqrt{c_i}$ values, 1 mark for n_i values].</p>	i	N_i	σ_i	c_i	$N_i\sigma_i/\sqrt{c_i}$	1	1000	1550	1.5	1265570	2	1100	2150	1.5	1931014	3	1600	2950	1.5	3853864	4	600	1100	1	660000	Total	4300			7710448
i	N_i	σ_i	c_i	$N_i\sigma_i/\sqrt{c_i}$																											
1	1000	1550	1.5	1265570																											
2	1100	2150	1.5	1931014																											
3	1600	2950	1.5	3853864																											
4	600	1100	1	660000																											
Total	4300			7710448																											
(ii)	$Var(\bar{x}_{prop}) = \sum \left(\frac{N_i}{N}\right)^2 \frac{\sigma_i^2}{n_i} = 10297.86 \quad [1 \text{ mark formula, 1 mark value}]$ $Var(\bar{x}_{opt}) = \sum \left(\frac{N_i}{N}\right)^2 \frac{\sigma_i^2}{n_i} = 9372.01 \quad [1 \text{ mark value}]$ $Var(\bar{x}_{srs}) = \frac{\sigma^2}{n} = \frac{3750^2}{500} = 28125 \quad [1 \text{ mark formula, 1 mark value}]$ $eff(\bar{x}_{prop}) = \frac{28125}{10297.86} = 2.73 \quad [1 \text{ mark formula, 1 mark value}]$ $eff(\bar{x}_{opt}) = \frac{28125}{9372.01} = 3.00 \quad [1 \text{ mark value}]$ <p style="margin-top: 10px;">Prefer optimal allocation which has the highest relative efficiency [1 mark].</p>																														
(iii)	<p>Without knowing the N and N_i's of interest in advance, the correct sample sizes cannot be calculated for either stratified sampling method [1 mark]. Using the wrong weights would introduce bias into the estimation of the overall mean income [1 mark]. It might be better to use simple random sampling which is unbiased [1 mark].</p>																														
(iv)	<p>It seems likely that graduate level jobs will have systematically higher salaries, so a further improvement in relative efficiency could be made by taking this categorisation into account [1 mark]. If the proportion of graduate-level jobs in each stratum were known in advance, then post-hoc stratification could be applied to the data collected [1 mark]. However, it seems that this information is not available and it would be dangerous to assume that the sample proportions would be good estimates of the population proportions, due to potential self-selection bias, so this further stratification cannot sensibly be carried out [1 mark].</p>																														