# SOLUTIONS

| Qn | Solution | Total |
|---|---|---|
| 1a | 210 records | |
| 1b | Data mining task | |
| 1c | (i)  14.8<br>(ii)  Canadian group mean area: 11.9<br>Kama group mean area: 14.3<br>Rosa group mean area: 18.3<br>(iii)  As we see overlaps over the range (Canadian and Rosa seen to have little overlap compared to Canadian and Kama on the same scale from the boxplot), we can conclude that area alone cannot accurately be used to distinguish between different seed species.<br><br><br><br>(iv)  Yes, because the two clusters are distinctly separated as shown in the scatterplot for 'Canadian' (triangle) and 'Rosa' (circle). 'Canadian' generally have small area and length of kernel groove; while 'Rosa' generally have large area and length of kernel groove.<br><br> | 10 |
| 2a | $k = 3$ because there are three different species of seeds. | |

# SOLUTIONS

| | | |
|---|---|---|
| 2b | Clustering was done on normalized data so as not to let any attribute dominate the distance calculation due to large numerical values. | |
| 2c | Number of records in cluster_0 = 63<br>Number of records in cluster_1 = 68<br>Number of records in cluster_2 = 76 | |
| 2d |  | |
| 2e | cluster_1<br>cluster_0<br>cluster_2 | 10 |
| 3a | Group | |
| 3b | (i)     Node: File Reader<br>(ii)    Node: Partitioning<br>(iii)   Node: Decision Tree Learner<br>(iv)   Node: Decision Tree Predictor<br>(v)    Node: Scorer | |
| 3c | <br>False Negative for 'Kama' = 2<br>It means that 2 records of Kama seed have been misclassified as Canadian seed by the decision tree model. | |

For 2d:

| Row ID | Area | Perimeter | Compa... | Length | Width | Asymm... | Length ... |
|---|---|---|---|---|---|---|---|
| cluster_0 | 0.76 | 0.796 | 0.697 | 0.732 | 0.773 | 0.365 | 0.759 |
| cluster_1 | 0.383 | 0.419 | 0.673 | 0.363 | 0.469 | 0.261 | 0.317 |
| cluster_2 | 0.123 | 0.174 | 0.377 | 0.186 | 0.161 | 0.499 | 0.279 |

For 3c:

| Row ID | Kama | Rosa | Canadian |
|---|---|---|---|
| Kama | 10 | 0 | 2 |
| Rosa | 0 | 17 | 0 |
| Canadian | 0 | 0 | 13 |

# SOLUTIONS

3d

Kama (58/168)

| Category | % | n |
|---|---|---|
| Kama | 34.5 | 58 |
| Rosa | 31.5 | 53 |
| Canadian | 33.9 | 57 |
| Total | 100.0 | 168 |

*Length of kernel gr...*

<= 5.5755                    > 5.5755

Kama (57/116)

| Category | % | n |
|---|---|---|
| Kama | 49.1 | 57 |
| Rosa | 1.7 | 2 |
| Canadian | 49.1 | 57 |
| Total | 69.0 | 116 |

Rosa (51/52)

| Category | % | n |
|---|---|---|
| Kama | 1.9 | 1 |
| Rosa | 98.1 | 51 |
| Canadian | 0.0 | 0 |
| Total | 31.0 | 52 |

*Area*

<= 12.71                    > 12.71

Canadian (50/54)

| Category | % | n |
|---|---|---|
| Kama | 7.4 | 4 |
| Rosa | 0.0 | 0 |
| Canadian | 92.6 | 50 |
| Total | 32.1 | 54 |

Kama (53/62)

| Category | % | n |
|---|---|---|
| Kama | 85.5 | 53 |
| Rosa | 3.2 | 2 |
| Canadian | 11.3 | 7 |
| Total | 36.9 | 62 |

*Asymmetry coeffic...*

<= 4.4135                    > 4.4135

Kama (49/52)

| Category | % | n |
|---|---|---|
| Kama | 94.2 | 49 |
| Rosa | 3.8 | 2 |
| Canadian | 1.9 | 1 |
| Total | 31.0 | 52 |

Canadian (6/10)

| Category | % | n |
|---|---|---|
| Kama | 40.0 | 4 |
| Rosa | 0.0 | 0 |
| Canadian | 60.0 | 6 |
| Total | 6.0 | 10 |

IF Length of kernel groove > 5.5755 THEN Group = 'Rosa' with probability 0.981.
IF Length of kernel groove <= 5.5755 AND Area <= 12.71 THEN Group = 'Canadian' with probability 0.926.
IF Length of kernel groove <= 5.5755 AND Area > 12.71 AND Asummetry coefficient <= 4.4135 THEN Group = 'Kama' with probability 0.942.

15

# SOLUTIONS

| | | |
|---|---|---|
| 4a | Predictor variable:   Systolic<br>Response varibale:  Diastolic | |
| 4b | **Statistics on Linear Regression**<br><br>| Variable | Coeff. | Std. Err. | t-value | P>\|t\| |<br>|---|---|---|---|---|<br>| Systolic | 0.7692 | 0.2542 | 3.026 | 0.0105 |<br>| Intercept | -14.3798 | 34.1176 | -0.4215 | 0.6809 |<br><br>Multiple R-Squared: 0.4328<br>Adjusted R-Squared: 0.3855<br><br>Estimated_Diastolic = 0.7692×Systolic − 14.3798 | |
| 4c | For every 1 mmHg increase in systolic pressure, the diastolic pressure is predicted to increase by 0.7692 mmHg.<br>The intercept value carries no meaning in this context. | |
| 4d | Estimated_Diastolic = 0.7692×120 − 14.3798 = 77.9 mmHg<br>(give full mark if only answer is given as students may use KNIME to generate the answer) | |
| 4e | $R^2$ is 0.4328 (not high) which suggests that the regression equation may not be useful for making estimation about the mean response variable. | 15 |
| 5a | (i)      6<br>(ii)      2<br>(iii)     Cluster_2, because distance to centroid_2 is smaller. | |
| 5b | (i)    Centroid of Cluster $1 = \dfrac{2+3+4}{3} = 3$<br><br>(ii)    Centroid of Cluster $2 = \dfrac{10+11+20+25+30+12}{6} = 18$ | |
| 5c | $k = 3$<br>This is because from the plot, the 'elbow' is indicated at $k$=3. It can be observed that by adding another cluster (say $k$=4) does not give much better modeling of the data in terms of reducing the WSSE. Hence, the best k is 3.<br>OR<br>This is because from the plot, the rate of decrease of WSSE sharply shifts when $k$=3 in the elbow plot. Hence, the best k is 3. | 10 |

# SOLUTIONS

**SINGAPORE POLYTECHNIC**
**2019 / 2020 Semester 1 EST**

Module Name: Statistics and Analytics for Engineers

Module Code: <u>MS_SAE</u>

Course:

Year: **2 FT**

**Page** 5 **of** 6

| 6a | $P(<=50K)=16/20$ | $P(>50K)=4/20$ | |
|---|---|---|---|
| | $P(Education \leq 9 \vert \leq 50K)=7/16$ | $P(Education \leq 9 \vert >50K)=1/4$ | |
| | $P(Education > 9 \vert \leq 50K)=9/16$ | $P(Education > 9 \vert >50K)=3/4$ | |
| | $P(Marital\_Status = Married \vert \leq 50K)=8/16$ | $P(Marital\_Status = Married \vert >50K)=2/4$ | |
| | $P(Marital\_Status = Others \vert \leq 50K)=8/16$ | $P(Marital\_Status = Others \vert >50K)=2/4$ | |
| 6b | Let $\mathbf{x} = (Education > 9\,years, Marital\_Status = Others)$ $$P(<=50K \vert \mathbf{x}) = \frac{P(\mathbf{x} \vert <=50K).P(<=50K)}{P(\mathbf{x})}$$ $$= \frac{P(Education>9 \vert <=50K).P(Marital\_Status=Others \vert <=50K).P(<=50K)}{P(\mathbf{x})}$$ $$= \frac{\left(\frac{9}{16}\right)\left(\frac{8}{16}\right)\left(\frac{16}{20}\right)}{P(\mathbf{x})} \approx \frac{0.225}{P(\mathbf{x})}$$ $$P(>50K \vert \mathbf{x}) = \frac{P(\mathbf{x} \vert >50K).P(>50K)}{P(\mathbf{x})}$$ $$= \frac{P(Education>9 \vert >50K).P(Marital\_Status=Others \vert >50K).P(>50K)}{P(\mathbf{x})}$$ $$= \frac{\left(\frac{3}{4}\right)\left(\frac{2}{4}\right)\left(\frac{4}{20}\right)}{P(\mathbf{x})} \approx \frac{0.075}{P(\mathbf{x})}$$ Since $\frac{0.225}{P(\mathbf{x})} > \frac{0.075}{P(\mathbf{x})}$, $\mathbf{x}$ is likely to earn <=50K per year. | | |
| 6c | With independence assumption: $$P(>9\,years \vert >50K)\,P(Others \vert >50K) = \frac{3}{4}.\frac{2}{4} = \frac{3}{8}$$ | | 24 |

# SOLUTIONS

| | | | |
|---|---|---|---|
| 7ai | $Y = -1312.2922 + 3.8165(X1) + 0.0444(X2) + 0.1525(X3) + 259.4698(X4)$ | 3 | |
| 7a ii | The most influential variable is X4 as its coefficent is considerably larger than those of other predictor variables. In addition, its p-value also indicates that it has statistical significance. | 2 | |
| 7a iii | The numerical meaning of the the coefficient of X4 implies that if 'Market Share' changes by a unit measure, the 'Sales' is expected to increase by 259.47 units. | 2 | |
| 7a iv | $Y = -1312.2922 + 3.8165(40) + 0.0444(50000) + 0.1525(6000)$ $+ 259.4698(12)$ $= 5089.0054 \text{ or } = 5090 \text{ (corr. to 3 s.f.)}$ | 3 | |
| 7a v | The $R^2$ value is 0.896. This means that the linear regression model is able to explain 89.6% of the variation observed in the data. (It is unable to explain 10.4% of the variation in the data, which is thus attributed to error-based, random variation.) | 2 | |
| 7f | The adjusted $R^2$ value adjusts the $R^2$ value for the number of predictor variables in the model. Its value increases only when a newly included predictor variable improves the model fit to the data by an amount more than what pure chance would do. So, the adjusted $R^2$ value may decrease when a newly included predictor variable does not improve the model fit by a sufficient amount. | 2 | |
| 7g | Introduction of X5 in the model has increased both $R^2$ (0.915) and adjusted $R^2$ values (0.8926). In addition, the size of the coefficient of X5 is the largest of all the coefficients. So, X5 has large impact on the expected prediction of 'Sales' and this impact is also significant as suggested by the increase in $R^2$ and adjusted $R^2$ values. | 2 | 16 |