# SOLUTIONS

**SINGAPORE POLYTECHNIC**
**2020 / 2021 Semester 2 EST**

Module Name: Statistics and Analytics for Engineers
Course:  DARE  DASE  DBEN  DCPE  DEB  DEEE  DME DMRO

Module Code: MS_SAE
Year: 2 FT
Page 1 of 2

**MCQ**

| 1. | 2. | 3. | 4. | 5. |
|---|---|---|---|---|
| c | b | b | d | b |

| Q1 | | Solution |
|---|---|---|
| a(i) | Number of rows — 1000 | |
| | Number of columns — 10 | |
| a(ii) | 2 | |
| b(i) | 5 | |
| b(ii) | gender, last_new_job | |
| c(i) | (A) 47 | |
| | (B) 60.171 | |
| | *Accept (B) 60.2* | |
| c(ii) | (A) 46 | |
| | (B) 64.882 | |
| | *Accept (B) 64.9* | |
| c(iii) | 22 | |
| d | Number of records after Step I — 812 | |
| | Number of records in 'Very Experienced' — 116 | |
| | % of samples in 'Very Experienced' — 14.29 % | |
| | *Accept (B) 14.3%* | |

| Q2 | | Solution | | |
|---|---|---|---|---|
| a | k=2, since the data consists of only 2 groups (G and F) | | | |
| b | The notes are relatively well separated using measurements on variance and skewness. Genuine notes on average appear to have higher values on these two measurements than fake notes. Measurements on entropy, on the other hand, do not appear to separate the notes well, as there is considerable overlap in their values between the two groups. | | | |
| c | Different variables may be measured on different scales. Normalization is important so as not to let any variable with large numerical values to dominate the distance calculations. | | | |
| d | | (I) | (II) | |
| | (A) | 7.502 | 7.50 | |
| | (B) | 0.696 | 0.70 | |
| | (C) | -9.583 | -9.58 | |
| | (D) | 0 | 0 | |
| | *Accept either (I) or (II)* | | | |

e

$$G = \left( \frac{0.696 + 0.909 + 1 + 0.982}{4} , \frac{0.739 + 1 + 0.721 + 0.36}{4} \right)$$
$$= (0.897, 0.705)$$

f

| | | |
|---|---|---|
| Standardized Variance | $\dfrac{2.20 - (-3.59)}{7.502}$ | 0.772 |
| Standardized Skewness | $\dfrac{6.00 - (-9.583)}{17.421}$ | 0.894 |
| $d(\text{newnote}, \text{centroid}_G)$ | $\sqrt{(0.772 - 0.897)^2 + (0.894 - 0.705)^2}$ | 0.227 |
| $d(\text{newnote}, \text{centroid}_F)$ | $\sqrt{(0.772 - 0.208)^2 + (0.894 - 0.303)^2}$ | 0.817 |

Since $d(\text{newnote}, \text{centroid}_G) < d(\text{newnote}, \text{centroid}_F)$ , the new note is likely to be genuine.
Note: Accept solution without standardisation.

# SOLUTIONS

**SINGAPORE POLYTECHNIC**
**2020 / 2021 Semester 2 EST**

Module Name: Statistics and Analytics for Engineers

Course: DARE DASE DBEN DCPE DEB DEEE DME DMRO

Module Code: MS_SAE

Year: 2 FT

Page 2 of 2

| Q3 | Solution |
|---|---|
| a | bad_loans |
| b | Target variable is categorical |
| c | This is a classification task where the target 'bad loans' is a class label, that is categorical. |
| d | Using 80% of the data (i.e., 8000 records) to train the model, and 20% of the remaining data (i.e., 2000 records) to evaluate the accuracy of the model. |
| e | Gini index |
| f | First split occurs with the attribute "term (in months) <=48" or "term (in months) >48". Splitting with this attribute gives the lowest Gini index (or greatest gain) compared to other attributes (with reference to the unsplitted records.) |
| g | 79.3% or 0.793 $\frac{1552+35}{2000} = 0.7935$ |
| h(i) | 68 68 borrowers were predicted to default payment of their loan but they did not. |
| h(ii) | 345 345 borrowers who were predicted not to default payment of their loan, defaulted. |
| h(iii) | Having a high false negative error might result in the United Finance losing potential lender customers as the lenders might end up with bad debts due to default of payment by borrowers who were predicted not to default their payment. *Accept other answers deemed logical* |
| | |

| Q4 | Solution | |
|---|---|---|
| a | House Price | |
| b(i) | If there is a unit increase in bedroom, the house price will drop by $62,829 | |
| b(ii) | Estimated Price = -120018 + 45392*floors + 319.94*sqft_living - 35.1*sqft_above – 62829*bedrooms + 57093*condition - 0.672*sqft_lot | |
| b(iii) | No. Because the R-sq and R-sq(adj) values are close. | |
| c | None all coefficients are < 5% | |
| d(i) | 48.45% | |
| d(ii) | 48.45% of the house price variation that is explained by the regression model | |
| e | sqft_lot*condition, $p > 0.05$ | |
| f | Steps | Name of the KNIME Node |
| | Create the interaction term | Math Formula |
| | To evaluate and score the model | Numeric Scorer |
| | | |