# Cover sheet for submission of work for assessment

SWiN BUR NE *

SWINBURNE UNIVERSITY OF TECHNOLOGY

## UNIT DETAILS

| | | | | | |
|---|---|---|---|---|---|
| Unit name | Data Science Principles | | Class day/time | Wed, 8 – 12am | Office use only |
| Unit code | COS10022 | Assignment no. 02 | Due date | 26/03/2023 | |
| Name of lecturer/teacher | Dr. Pham Thi Kim Dung | | | | |
| Tutor/marker's name | Dr. Pham Thi Kim Dung | | | | Faculty or school date stamp |

## STUDENT(S)

| Family Name | Given Name | Student ID Number |
|---|---|---|
| Hau | Linh Chi | 104177160 |

## DECLARATION AND STATEMENT OF AUTHORSHIP

1. For the sake of this evaluation, I have not impersonated anyone or let anyone else to impersonate me.
2. This evaluation is all original work from myself, with the exception of the places where proper credit has been given.
3. Except where such collaboration has been approved by the lecturer or instructor in question, no portion of this evaluation has been prepared for me by anyone else.
4. I have not previously submitted this work for this or any other course/unit.
5. I accept that my assessment response may be duplicated, shared, compared, stored, and used for benchmarking, plagiarism detection, or educational purposes.

I understand that:

6. Plagiarism is known as the act of presenting another person's work, idea, or creativity as your own. It is a sort of academic fraud that might get you kicked out of the university and is considered cheating. Written, graphic, and visual works, computer data, oral presentations, and other forms of presentation can all be used to source and deliver plagiarized content. When the source of the copied material is not properly cited, plagiarism occurs.

**Student signature/s**

I declare that I have read and understood the declaration and statement of authorship.

COS10022 – Data Science Principles – Assignment 1

# DATA CLEANING AND ANALYTICS

## I. ASSIGNMENT SUMMARY

This assignment provides an overview of data cleaning and prediction model construction, including the key concepts, procedures, and tools involved.
The focus is on:
- Selecting appropriate features and models while gaining some knowledge about prebuilt tools and applying them in a data science project.
- Preparing the dataset using KNIME analytical platform for cleaning and two classification models, Naïve Bayes and Random Forest, are developed.
- Selecting relevant attributes, cleaning the dataset, partitioning the data into training and test sets, developing an effective prediction model, and explaining the results.

The goal of the assignment is to gain hands-on experience with data cleaning and model construction in a real-world setting.

## II. INTRODUCTION

The focus of this report is a dataset that was gathered from the real world, consisting of 100,000 tuples that are categorised into three different financial credit score classes. The original data includes 24 attributes in total.
The goals of this assignment are: firstly, to carry out the necessary data cleaning and preparation for future use, and secondly, to develop two predictive models that can be used to project the "Credit_Score" classification.
My works are described in detail in this report and include:
- Preparing the raw data for future use by cleaning and organising it.
- Constructing two models (Naïve Bayes and Random Forest classifier models) to forecast the value and categorise financial credit scores.
- Choosing and implementing suitable features and models for this data project.
- Identifying relevant attributes, splitting the dataset into training and testing sets, developing predictive models, and explaining the results.

## III. DATA CLEANING

Question 1.1

| Node | Configuration |
|------|---------------|

| | | |
|---|---|---|
| **Column Filter**<br> | The excluded attribute is *"Name"*. The reason for removing it are:<br>• Credit scores are determined by a variety of factors such as payment history, credit utilisation, length of credit history, types of credit used, new credit inquiries, etc. Therefore, a person's name does not affect their credit score.<br>• In the Naïve Bayes learner process, the maximum number of unique nominal values per attribute is set to 600 to limited by the amount of available data, the complexity of the problem, and the performance requirements of the application. However, the number of names in the dataset exceeds the limit number (10137 different names). So, the *"Name"* attribute is filtered out. |  |



## Question 1.2

| S | Node | Configuration |
|---|---|---|
| 1 | **Missing Value**<br> | This node is used to remove tuples with missing values from these 18 attributes *"Month"*, *"Age"*, *"Occupation"*, *"Annual_Income"*, *"Num_Bank_Accounts"*, *"Num_Credit_Card"*, *"Interest_Rate"*, *"Num_of_Loan"*, *"Delay_from_due_date"*, *"Changed_Credit_Limit"*, *"Credit_Mix"*, *"Outstanding_debt"*, *"Credit_Utilization_Ratio"*, *"Credit_History_Age"*, *"Payment_of_Min_Amount"*, *"Total_EMI_per_month"*, *"Amount_invested_monthly"* and *"Payment_Behaviour"*.<br><br>One example output of the *"Amount_invested_monthly" column*: |  |

| S | Node | Configuration |
|---|---|---|
| | |  |
| 2 | **Rule-based Row Filter**  | Some tuples containing infeasible values, such as: <br> • *"Monthly_Inhand_Salary"* < 0 <br> • *"Num_Bank_Accounts"* < 0 <br> • *"Num_Credit_Card"* < 0 <br> • *"Changed_Credit_Limit"* contains "_" <br> are removed by using this node with the commands (exclude TRUE matches) as follows: <br> $Monthly_Inhand_Salary$ < 0 => TRUE <br> $Num_Bank_Accounts$ < 0 => TRUE <br> $Num_Credit_Card$ < 0 => TRUE <br> $Changed_Credit_Limit$ MATCHES "_" => TRUE <br><br> Example output of the *"Changed_Credit_Limit"* column:  <br>  |

## Question 1.3

| S | Node | Configuration |
|---|---|---|
| 1 | **String Manipulation**  | Symbols that are not numbers are eliminated from the *"Age"* attribute and the data of this attribute is converted into the usual number format (to integer) and replace the new output to the original *"Age"* column. The "." and "-" are kept because they represent decimal number and negative number (some of them must be remove, so erasing the "-" can lead to wrong data cleaning, respectively. <br> Expression: <br> toInt(regexReplace($Age$,"[^0-9.-]","")) <br><br> Example output: |

| Node | Configuration |
|---|---|
| **2**<br><br>**Rule-based Row Filter** | The tuples having *"Age"* value lower than or equal to 0 or greater than 120 are dropped using the following command (include TRUE matches):<br>$Age$ > 0 AND $Age$ <= 120 => TRUE<br><br>Example output:<br> |

## Question 1.4

| Node | Configuration |
|---|---|
| **String Manipulation** | Non-numerical symbols in the *"Annual_Income"* column are eliminated and the remaining data are converted to the double format. The "." and "-" are kept because they represent decimal number and negative number, respectively. The new output is replaced to the original *"Annual_Income"* column.<br>Expression:<br>toDouble(regexReplace($Annual_Income$,"[^0-9.-]","")) <br><br>Example output:<br> |

Question 1.5

| S | Node | Configuration |
|---|------|---------------|
| 1 | **String Manipulation** r[s] | "_____" values from the *"Occupation"* attribute are changed to null. The original *"Occupation"* column is replaced with the new data. Expression: replace($Occupation$,"_____", null)  |
| 2 | **String Manipulation** r[s] | Non-numerical symbols in *"Num_of_Loan"* are removed and the data is converted to integer data type. The "." and "-" are kept because they represent decimal number and negative number, respectively. The new output is replaced to the original *"Num_of_Loan"* column. Expression: toInt(regexReplace($Num_of_Loan$,"[^0-9.-]",""))  Example output:  |
| 3 | **Math Formula (Multi Column)** f(x) | This node is utilised to take the absolute values of attributes *"Num_Bank_Accounts"* and *"Num_Credit_Card"* using the command as follows: abs($$CURRENT_COLUMN$$) The original *"Num_Bank_Accounts"* and *"Num_Credit_Card"* columns are replaced with the new data. (negative values of those two attributes are removed by the Rule-based Row Filter node in the question 2)  |

| | | |
|---|---|---|
| 4 | **Math Formula**<br><br>f(x) | For the *"Num_of_Loan"* attribute, if the original values are negative, they are set to 0 and the new data is replaced to the original *"Num_of_Loan"* column.<br>Expression:<br>if($Num_of_Loan$<0,0, $Num_of_Loan$)<br><br>Example output:<br><br>**File Table - 3:1 - CSV Reader**<br>File  Edit  Hilite  Navigation  View<br>Table "default" - Rows: 100000  Spec - Colur<br><br>| Row ID | S ▾ Num_of_Loan |<br>|---|---|<br>| Row27 | 1 |<br>| Row28 | 1 |<br>| Row29 | 1 |<br>| Row30 | 1 |<br>| Row31 | -100 |<br><br>**Output data - 3:11 - Math Formula**<br>File  Edit  Hilite  Navigation  View<br>Table "default" - Rows: 90928  Spec - Colur<br><br>| Row ID | I Num_of_Loan |<br>|---|---|<br>| Row28 | 1 |<br>| Row29 | 1 |<br>| Row30 | 1 |<br>| Row31 | 0 | |
| 5 | **String Manipulation**<br><br>r[s] | Non-numerical symbols in *"Num_of_Delayed_payment"* are taken out and the data is converted to integer.  The "." and "-" are kept because they represent decimal number and negative number, respectively. The new output is replaced to the original *"Num_of_Delayed_payment"* column.<br>Expression:<br>toInt(regexReplace Num_of_Delayed_Payment$,"[^0-9.-]",""))<br><br>Example output:<br><br>**File Table - 3:1 - CSV Reader**<br>File  Edit  Hilite  Navigation  View<br>Table "default" - Rows: 100000  Spec - Columns: 25  P<br><br>| Row ID | S Num_of_Delayed_Payment |<br>|---|---|<br>| Row51 | 2 |<br>| Row52 | 4 |<br>| Row53 | 3_ |<br>| Row54 | 2_ |<br>| Row55 | 2 |<br><br>**Appended table - 3:12 - String Manipulation**<br>File  Edit  Hilite  Navigation  View<br>Table "default" - Rows: 90928  Spec - Columns: 24  P<br><br>| Row ID | I Num_of_Delayed_Payment |<br>|---|---|<br>| Row52 | 4 |<br>| Row53 | 3 |<br>| Row54 | 2 |<br>| Row55 | 2 |<br>| Row57 | 14 | |

| | | |
|---|---|---|
| 6 | **String Manipulation**<br>r[s] | The *"Credit_Mix"* values are set to "Unknown" if the original value is "_" using the following expression:<br>replace($Credit_Mix$, "_","Unknown")<br>The new output is replaced to the original *"Credit_Mix"* column.<br><br>Example output: |
| 7 | **String Manipulation**<br>r[s] | The non-numerical symbols in *"Outstanding_Debt"* attribute are erased and the data is converted to the double format. The "." and "-" are kept because they represent decimal number and negative number, respectively. The new output is replaced to the original *"Outstanding_Debt"* column.<br>Expression:<br>toDouble(regexReplace Outstanding_Debt$,"[^0-9.-]",""))<br><br>Example output: |





<span style="background-color:yellow">Question 1.6</span>

| S | Node | Configuration |
|---|---|---|

| | |
|---|---|
| **String Manipulation**<br>r[s] | This node is applied to separate the number of year from the *"Credit_History_Age"* string value, erase white spaces and convert it to the number format.<br>The separated output is appended to the new column named *"Years"*.<br>Expression:<br>toInt(strip(substr<br>($Credit_History_Age$,0,2)))<br>Two characters counts from the first character (index 0) of a *"Credit_History_Age"* value are taken by using the expression substr(). For example, "9 Years and 6 Months" returns "9 " after substruction. The strip() expression is utilised to remove white spaces. For example, "9 " returns "9" after this function.<br><br>Example output: |

| | |
|---|---|
| **String Manipulation**<br>r[s] | This node is applied to separate the number of month from the *"Credit_History_Age"* string value, erase white spaces and convert it to the number format.<br>The separated output is appended to the new column named *"Months"*.<br>Expression:<br>toInt(strip(substr<br>Credit_History_Age$,indexOfChars<br>($Credit_History_Age$,"d")+2,2)))<br>Two characters counts from the first character starting after two characters from the letter "d" (indexOfChars($Credit_History_Age$,"d")+2) of a *"Credit_History_Age"* value are taken by using the expression substr(). For example, "9 Years and 6 Months" returns "6 " after substruction. The strip() expression is utilised to remove white spaces. For example, "6 " returns "6" after this function.<br><br>Example output: |

| | Math Formula | The *"Credit_History_Age"* values are converted to the count of months and store it in the integer format using the following formula: $Years$ * 12 + $Months$ The converted output is stored in a new column called *"Total_CHA"*.  Example output:  |

## Question 1.7

| S | Node | Configuration |
|---|------|---------------|
| 1 | String Manipulation (Multi Column) |  This node is used to remove the non-numerical symbol in two columns: *"Amount_invested_monthly"* and *"Monthly_Balance"* and convert it to the double format. Expression: toDouble(regexReplace( $$CURRENTCOLUMN$$,"[^0-9.-]","")) This function is used to apply regex to string and replaces string if regex matches. |

The "." and "-" are kept because they represent decimal number and negative number, respectively. The new output is replaced to their two original columns.

Example output:



---

**2**

**String Replacer**

This node is used to set the value to "Unknown" if the original value in *"Payment_Behaviour"* column starts with "!@". The "*" is applied to find the wildcard pattern in the right order.



Example output:

| 3 | String Manipulation | *"Changed_Credit_Limit"* colum is converted to the double format using the expression as follows:<br>toDouble(<br>$Changed_Credit_Limit$)<br>The new output is replaced to the original *"Changed_Credit_Limit"* column. |  |

## Question 1.8

| S | Node | Configuration |
|---|------|---------------|
| 1 | Missing Value | This node is used to replace the missing values in all string type attributes to the "Next Value*" and all numeric format to the "Previous Value*".<br><br>Example output:<br> |  |
| 2 | Math Formula | The *"Monthly_Balance"* value is replaced with 0 if it is negative using the if condition formula if(condition, true value, false value):<br>if($Monthly_Balance$<0,<br>0,$Num_of_Loan$)<br>The new output is replaced to the original *"Monthly_Balance"* columns. |  |

<mark>Question 1.9</mark>

| S | Node | Configuration |
|---|------|---------------|
| 1 | **String Manipulation**<br>f[s] | The *"Type_of_Loan"* column is simplified using this node. If the original content has more than one type separated by a comma, keep only the first part. Otherwise, keep the full description if there is no comma included. The output result is appended to a new column called *"tmp_ToL"*.<br>Expression:<br>substr($Type_of_Loan$, 0, indexOf($Type_of_Loan$, ","))<br><br>This function is applied to substract the string from *"Type_of_Loan"* column starting from the first character (index 0) to the first comma. For example, "Auto Loan, Auto Loan, and Not Specified" returns "Auto Loan" after substraction and "Credit-Builder Loan" returns missing value because there are not any comma in the value.<br><br>Example output:<br> |
| 2 | **Rule Engine**<br>☑ | This node is used to recover the missing value from the previous node to its original value from the *"Type_of_Loan"* column and keep the value subtracted in the previous node.<br>The new output is replaced to the *"tmp_ToL"* column.<br>Expression:<br>$tmp_ToL$ LIKE "" => $Type_of_Loan$ (the missing values are recovered from the original *"Type_of_Loan"* attribute)<br>TRUE => $tmp_ToL$ (the default outcome remains the same)<br><br><br>Example output:<br> |

## Question 1.10

| Node | Configuration |
|------|---------------|
| **Numeric Binner** | This node is utilised to bin the *"Changed_Credit_Limit"* attribute with six bins of ranges: $[-\infty,-3.0)$, $[-3.0,0)$, $[0,3.0)$, $[3.0,6.0)$, $[6.0,7.5)$, and $[7.5,\infty)$ and put the result into a new attribute called *"Changed_Credit_Limit_binned"*.<br><br>Example output: |

## Question 1.11

| S | Node | Configuration |
|---|------|---------------|
| 1 | **Column Filter** | The Column Fiter node is used to filter out the original columns that have their new data being appended to new columns (*"Type_of_Loan"*, *"Changed_Credit_Limit"*, *"Credit_History_Age"*) and attributes that are the output of the number of year and month substraction being used for converting the *"Credit_History_Age"* to the count of months (*"Years"*, *"Months"*). |
| 2 | **Feature Selection Loop Start (1:1)** | This node is applied to let the algorithm select the attributes using the Genetic Algorithm feature selection strategy with default population size and the maximum number of generations. The class label *"Credit_Score"* is excluded and the static random seed of 3122 is applied.<br><br>Example output: |

| | | | |
|---|---|---|---|
| 3 | **Shuffle** | The dataset is shuffled with the random seed being set to 3122. | *Dialog - 3:29 - Shuffle* — □ ✕<br>File<br>Seed  Flow Variables  Job Manager Selection  Memory Policy<br>☑ Use seed<br>3122        Draw new seed<br>OK    Apply    Cancel |
| 4 | **Partitioning** | The dataset is partitioned by Linear sampling and a 75:25 ratio, with 75% belongs to the training set and the remaining 25% for the test set. | *Dialog - 3:30 - Partitioning* — □ ✕<br>File<br>First partition  Flow Variables  Job Manager Selection  Memory Policy<br>Choose size of first partition<br>○ Absolute                100<br>● Relative[%]               75<br>○ Take from top<br>● Linear sampling<br>○ Draw randomly<br>○ Stratified sampling     S Changed_Credit_Limit_binned<br>☐ Use random seed         1,679,645,615,:<br>OK    Apply    Cancel |

# IV. DATA ANALYTICS

## 1. Naïve Bayes Classifier

<mark>Question 2.1</mark>

A screenshot of the Naïve Bayes classifier in the KNIME workflow:



<mark>Question 2.2</mark>

The default probability value and the minimum standard deviation is set to 0.0001. The threshold standard deviation is 0 and the maximum number of unique nominal values per attribute of 600 is applied to the classifier leaner.

*Dialog - 3:31 - Naive Bayes Learner* — □ ✕
File
Options  Flow Variables  Job Manager Selection  Memory Policy
Classification Column: S Credit_Score
Default probability: 0.0001
Minimum standard deviation 0.0001
Threshold standard deviation 0.0
Maximum number of unique nominal values per attribute: 600
☐ Ignore missing values    ☐ Create PMML 4.2 compatible model
OK    Apply    Cancel

<mark>Question 2.3</mark>

Screenshots of the Confusion Matrix and the Accuracy statistics of the test result:

**Confusion Matrix - 3:32 -...** — □ ✕
File  Hilite

| Credit_Sco... | Good | Standard | Poor |
|---|---|---|---|
| Good | 3376 | 620 | 107 |
| Standard | 3069 | 6651 | 2363 |
| Poor | 1034 | 1636 | 3876 |

Correct classified: 13,903    Wrong classified: 8,829

Accuracy: 61.16%    Error: 38.84%

Cohen's kappa (κ): 0.404%

**Accuracy statistics - 3:32 - Scorer**
File  Edit  Hilite  Navigation  View
Table "default" - Rows: 4  Spec - Columns: 11  Properties  Flow Variables

| Row ID | I TruePo... | I FalsePo... | I TrueNe... | I FalseN... | D Recall | D Precision | D Sensitivity | D Specificity | D F-meas... | D Accuracy | D Cohen'... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Good | 3376 | 4103 | 14526 | 727 | 0.823 | 0.451 | 0.823 | 0.78 | 0.583 | ? | ? |
| Standard | 6651 | 2256 | 8393 | 5432 | 0.55 | 0.747 | 0.55 | 0.788 | 0.634 | ? | ? |
| Poor | 3876 | 2470 | 13716 | 2670 | 0.592 | 0.611 | 0.592 | 0.847 | 0.601 | ? | ? |
| Overall | ? | ? | ? | ? | ? | ? | ? | ? | ? | 0.612 | 0.404 |

This Naïve Bayes classifier produces a low precision result of 0.451 which indicates that this classifier <mark>performs unsatisfactorily.</mark>

## <mark>Question 2.4</mark>

The measurement being looked at to interpret the conclusion in this case is: <mark>Precision</mark>.

If the bank wants to minimise the risk of lending money to customers, the "Good" in *"Credit_Score"* should be the major target. So, the statistic represents the value being predicted as "Good" are actually "Good" is taken into consideration.

In a confusion matrix, precision is defined as the ratio of true positives (TP: The number of "Good" customers is correctly classified) to the sum of true positives and false positives (FP: The number of "Poor"/"Standard" customers is incorrectly identified as "Good" ones):

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

Precision measures the proportion of positive predictions (which means the instances that the model classifies as positive) that are actually correct.

A high precision value indicates that the model is making very few false positive predictions (a few number of "Poor"/"Standard" customer incorrectly classified as "Good" ones) and vice versa.

This Naïve Bayes classifier produces a low precision result of 0.451 which indicates that this classifier perform unsatisfactorily.

## 2. Random Forest Classifier

### <mark>Question 3.1</mark>

A screenshot of the Random Forest classifier in the KNIME workflow:



### <mark>Question 3.2</mark>



### <mark>Question 3.3</mark>

The measurement being looked at to interpret the conclusion in this case is: <mark>Precision</mark>.

If the bank wants to minimise the risk of lending money to customers, the "Good" in *"Credit_Score"* should be the major target. So, the statistic represents the value being predicted as "Good" are actually "Good" is taken into consideration.

In a confusion matrix, precision is defined as the ratio of true positives (TP: The number of "Good" customers is correctly classified) to the sum of true positives and false positives (FP: The number of "Poor"/"Standard" customers is incorrectly identified as "Good" ones):

$$Precision = \frac{True\ Positive}{True\ Positive\ +\ False\ Positive}$$

Precision measures the proportion of positive predictions (which means the instances that the model classifies as positive) that are actually correct.
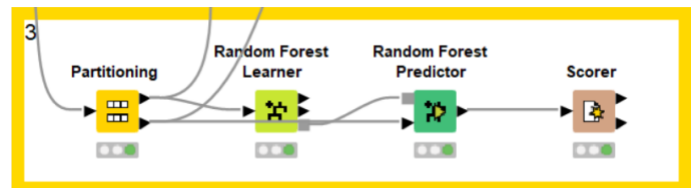
A high precision value indicates that the model is making very few false positive predictions (a few number of "Poor"/"Standard" customer incorrectly classified as "Good" ones) and vice versa.

Random Forest results and Naïve Bayes results Precision comparison:

|  | Random Forest | Naïve Bayes |
|---|---|---|
| Precision (Good) | 0.691 | 0.451 |
| Accuracy (Overall) | 0.752 | 0.612 |

The Naïve Bayes classifier produces a lower precision (and accuracy) result than the Random Forest one, which indicates that Random Forest model presents a more suitable result compared to Naïve Bayes one.

Question 3.4

Some measurements that should be looked at to find out which class performs the best by the model:

|  | Precision | Recall | F-measure |
|---|---|---|---|
| Good | 0.691 | 0.696 | 0.693 |
| Standard | 0.777 | 0.784 | 0.780 |
| Poor | 0.745 | 0.728 | 0.736 |

Precision, Recall and F-measure are common metrics used to evaluate the performance of classification models, including the Random Forest model:

- Precision is used to evaluate how well a model can correctly identify instances of a particular class. It measures the proportion of true positives (instances of the target class that were correctly identified by the model) out of all instances that the model classified as positive for that class. A higher precision means that the model makes fewer false positive errors, which means fewer instances that do not belong to the target class are mistakenly identified as belonging to it.
- Recall measures is used to evaluate how well a model can identify all relevant instances of a particular class. In other words, recall measures the proportion of true positives (instances of the target class that were correctly identified by the model) out of all actual positives (all instances of the target class in the dataset). In this classification problem where all classes are equally important, a higher recall for a particular class would indicate that the model is better at identifying that class.
- F-measure is a metric that combines precision and recall into a single score, which gives equal weight to both measures. It is calculated as the harmonic mean of precision and recall, with values ranging from 0 to 1, where a higher score indicates better performance. Therefore, a higher F-measure for a particular class means that the model is better at identifying instances of that class while also minimising false positives.

Looking at the results, it is obviously that <mark>class "Standard" performs the best results</mark>.


## V. CONCLUSION

In conclusion, this assignment aimed to clean and prepare raw data for future use, and to create two predictive models capable of projecting *"Credit_Score"* classifications. The tasks involved in this project have been thoroughly described in this report, including preparing the data by organising and cleaning it, constructing two models (Naïve Bayes and Random Forest classifier models) to predict and categorise financial credit scores, selecting appropriate features and models for the project, identifying relevant attributes, dividing the dataset into training and testing sets, building predictive models, and providing an explanation of the results. Overall, the report comprehensively showcases and explains the data cleaning and model constructing process.