**UTM**

UNIVERSITI TEKNOLOGI MALAYSIA

Faculty of Computing

UNIVERSITI PENYELIDIKAN

# UNIVERSITI TEKNOLOGI MALAYSIA
## SEMESTER II - 2020 / 2021

### ALTERNATIVE ASSESSMENT (INDIVIDUAL)

| | | |
|---|---|---|
| SUBJECT CODE | : | SCSP2753 |
| **SUBJECT NAME** | : | **DATA MINING** |
| **DATE** | : | **23 JUNE – 1 JULY 2021** |
| DURATION | : | 7 DAYS |
| SUBMISSION DATE | : | 1 JULY 2021 |

_____

INSTRUCTIONS:

This alternative assessment is assessed individually.
This alternative assessment consists of **TWO** parts.
Read **ALL** questions carefully and please answer **ALL** questions.
Submit  your answer (report) via e-learning and submission later than the due date is not accepted.
Any form of plagiarism is not allowed.

| | |
|---|---|
| Name | NUR ALEEYA SYAKILA BINTI MUHAMAD SUBIAN |
| I/C No. | 000726100810 |
| Year / Course | 2 SECP |
| Section | 01 |
| Lecturer Name | DR ROZILAWATI BINTI DOLLAH |

# Table Of Contents

# 1.0 PART 1

## 1.1 Introduction

There are a lot of tools out there that we can use to do Preprocessing and Data Mining Task. For this alternative assessment, I will be choosing RapidMiner as my tool. This is because RapidMiner has a lot of resources that I can refer to incase of I couldn't understand or misleading concept. This alternative assessment required me to do text preprocessing and few data mining tasks. To perform any data mining algorithm, preprocess the data is a must and it must be done before proceeding with other steps. In this alternative assessment, for question 1, I will be performing text preprocessing for appendix 1 given in the question sheet. After did some research about preprocessing, I found that every basic preprocessing has the same steps which are tokenize, filter stopwords, generate n gram, filter token and stemming. While for question 2, based on the appendix 2, the dataset that is given to me is dataset 29. Before performing the supervised and unsupervised learning algorithm, I preprocess my dataset beforehand.

# 1.1.1 Text Preprocessing

## Step 1 :

The data is in unstructured form where its format is variety such as in long text email or else. Therefore, I needed to convert it to the structured one by showing text only. To make it happened, I will be needed to do text preprocessing in RapidMiner. Firstly, I need to change the dataset from docx file into a xlsx file to make it easier for the RapidMiner to read the dataset.

## Input :

**APPENDIX 1**

**Protection of the small intestine from nonocclusive mesenteric ischemic injury due to cardiogenic shock.**

In a pericardial tamponade model of cardiogenic shock in pigs, we had previously shown that acute reductions in cardiac output produce severe mesenteric ischemia due to disproportionate splanchnic vasoconstriction. In this study, we extended the period of cardiogenic shock in order to investigate the pathogenesis of ischemic injury to the small intestinal wall. Four hours of tamponade produced sustained changes in splanchnic hemodynamics, similar to those observed in the prior short-term experiments. The resultant mesenteric ischemia caused necrotic lesions of the small intestine which were characteristic of those seen in nonocclusive mesenteric ischemia in human subjects. Prior alpha-adrenergic blockade failed to prevent either sustained mesenteric vasospasm or ischemic injury. In contrast, prior blockade of the renin-angiotensin axis, whether by nephrectomy or angiotensin-converting enzyme inhibition, blocked the splanchnic vasoconstriction, and thereby protected the small intestine from ischemic injury. The primary hemodynamic and pathologic features of this model of nonocclusive mesenteric ischemia appear to be mediated by the renin-angiotensin axis.

Figure 1 : Raw dataset appendix 1 in docx

## Output :

1. Protection of the small intestine from nonocclusive mesenteric ischemic injury due to cardiogenic shock.In a pericardial tamponade model of cardiogenic shock in pigs, we had previously shown that acute
2. reductions in cardiac output produce severe mesenteric ischemia due to disproportionate splanchnic vasoconstriction. In this study, we extended the period of cardiogenic shock in order to investigate the
3. pathogenesis of ischemic injury to the small intestinal wall. Four hours of tamponade produced sustained changes in splanchnic hemodynamics, similar to those observed in the prior short-term experiments.
4. The resultant mesenteric ischemia caused necrotic lesions of the small intestine which were characteristic of those seen in nonocclusive mesenteric ischemia in human subjects. Prior alpha-adrenergic blockade
5. failed to prevent either sustained mesenteric vasospasm or ischemic injury. In contrast, prior blockade of the renin-angiotensin axis, whether by nephrectomy or angiotensin-converting enzyme inhibition, blocked the splanchnic vasoconstriction, and thereby protected the small intestine from ischemic injury. The primary hemodynamic and pathologic features of this model of nonocclusive mesenteric ischemia appear to be mediated by the renin-angiotensin axis.

Figure 2 : Raw dataset appendix 1 in xlsx

## Step 2 : Read Excel

After converting it to a data in xlsx file, the very first basic thing I need to do after opening the RapidMiner is to drag the 'Read Excel' from the Operators as the dataset is in excel file. Drag the Read Excel operator and read the dataset file.
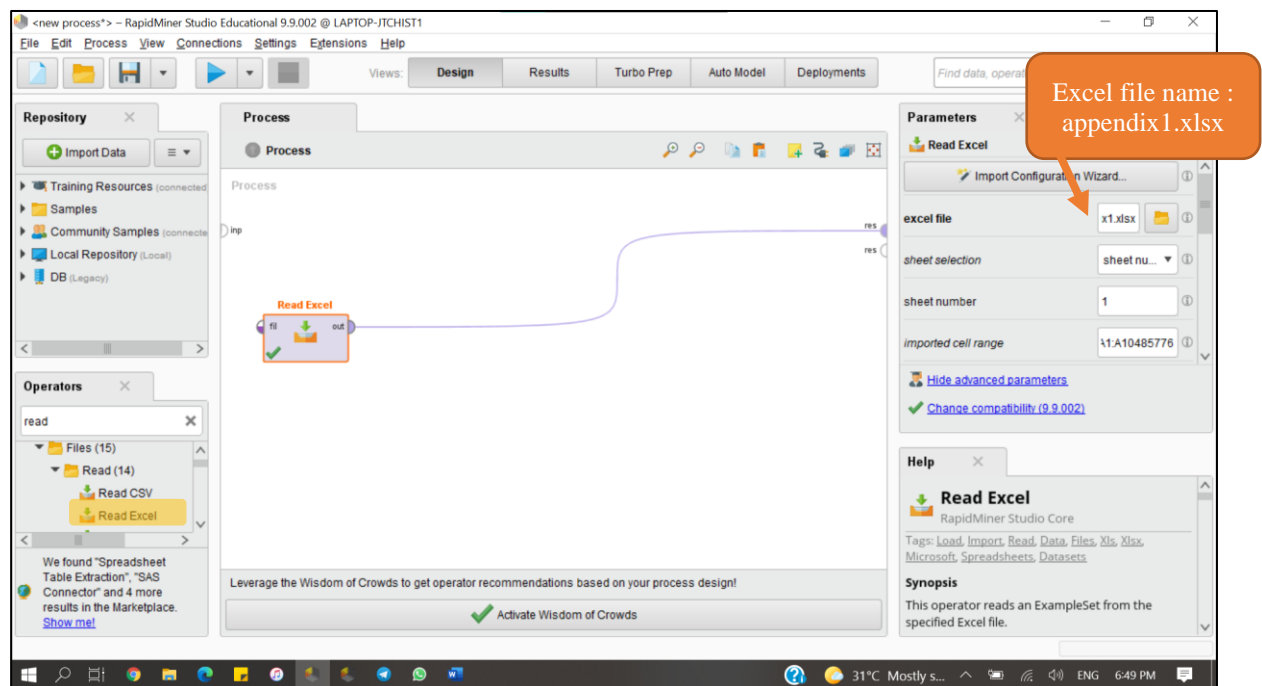
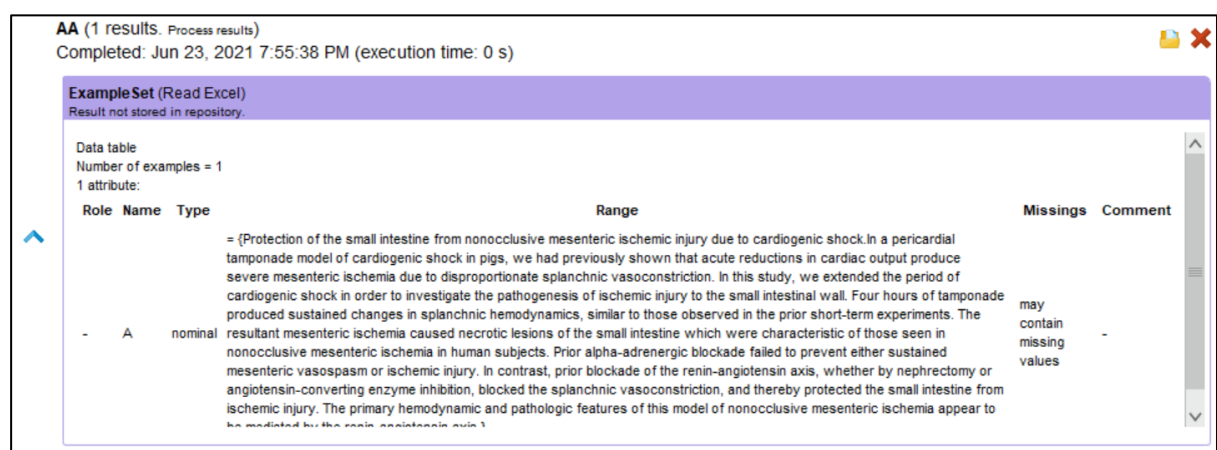## Input :



Figure 3 : Choose Read Excel operator

## Output :



Figure 4 : The result after Read Excel

## Step 3 : Nominal to Text

Then I needed to drag the 'Nominal to Text' operators. The Nominal to Text operator converts all nominal attributes to string attributes. Drag the 'Nominal to Text' operator. Connect the *output* of the 'Read Excel' operator to the *example* of the 'Nominal to Text' operator. Connect the output of 'Read Excel' operator to the example (input) of the 'Nominal to Text' operator

## Input :
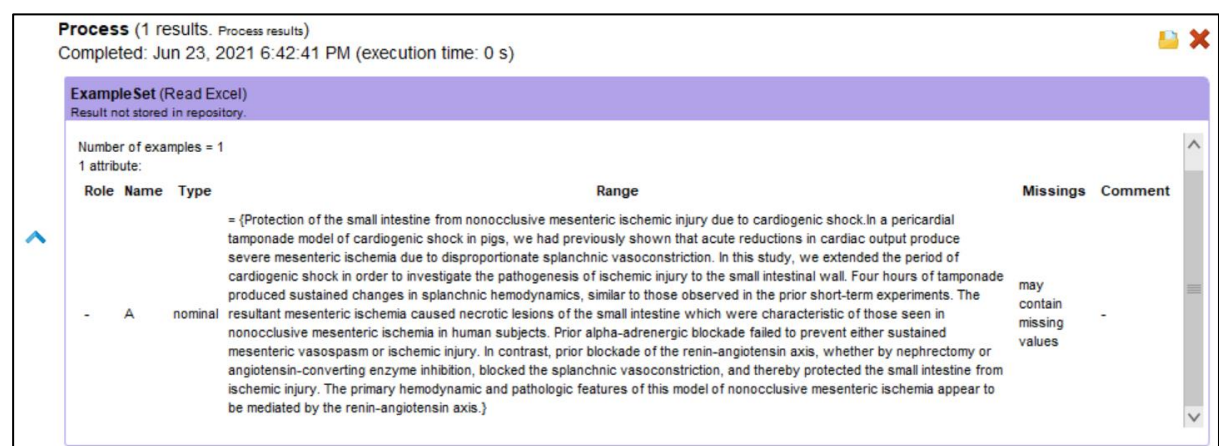


Figure 5 : Choose the Nominal to Text operator

## Output :



Figure 6 : The result after running the Nominal to Text operator

**Step 4 : Tokenize**

To continue the preprocessing, I needed to add 'Process Documents from Data' operator first. After that, proceed to the tokenize process. Tokenization is a pre-processing strategy what breaks a stream of text into words, expressions, images, or other significant components called tokens. After clicking the 'Process Document from Data', drag the 'Tokenize' operator to proceed.
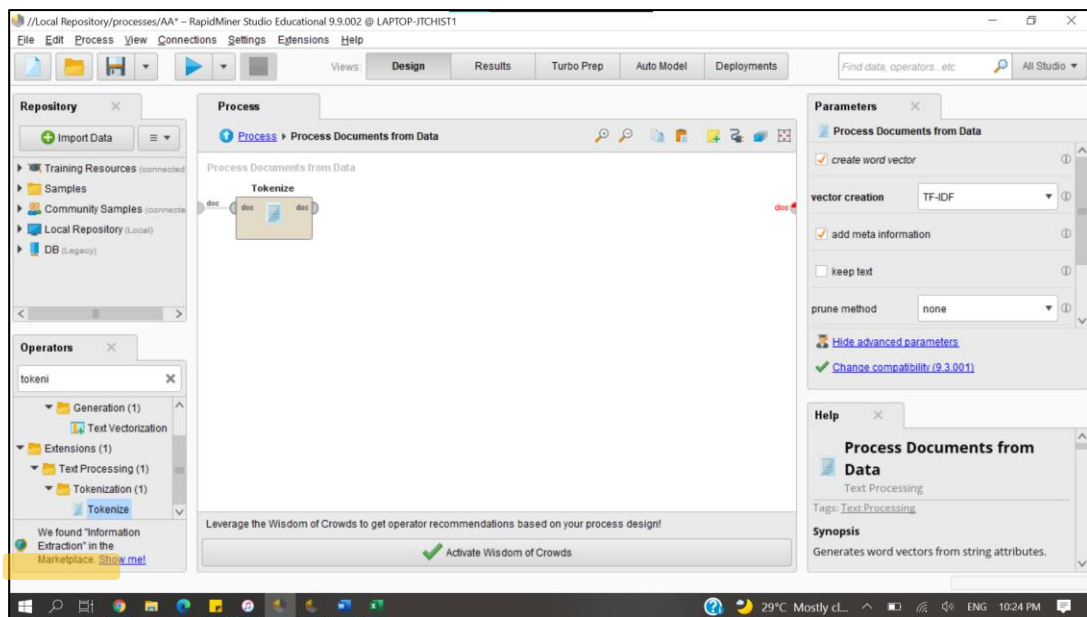
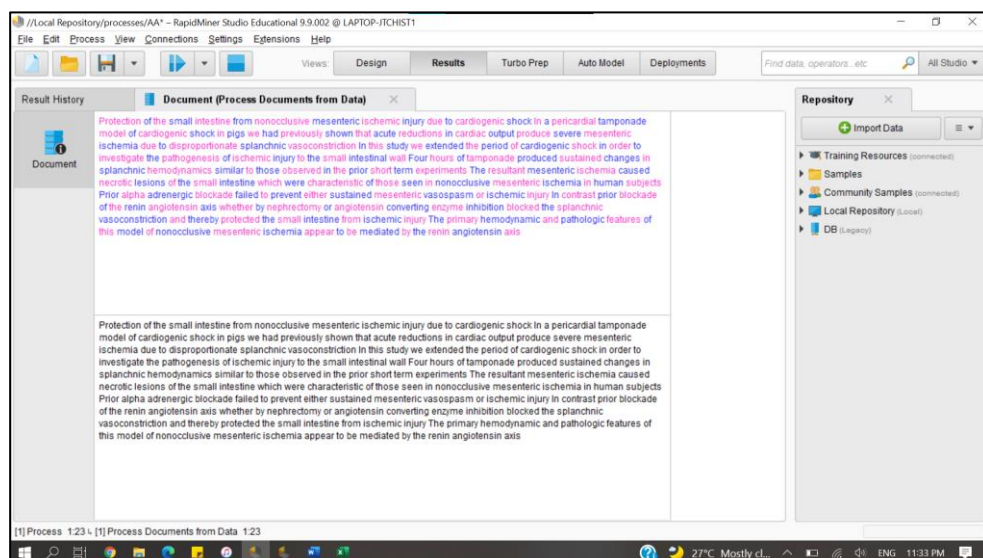**Input** :



Figure 8 : Add 'Tokenize' operator

**Output** :



Figure 9 : The result after the execution of the tokenize step

## Step 5 : Lower Cases

Then I needed to transform the dataset into lower cases. This step can help in situations where the dataset isn't exceptionally huge and fundamentally assists with consistency of anticipated output. Drag the Transform Cases operator into the process and make sure it is set to transform to lower case. After the execution, a list of lower-case words has been created. There are total of 96 words produced in the list.
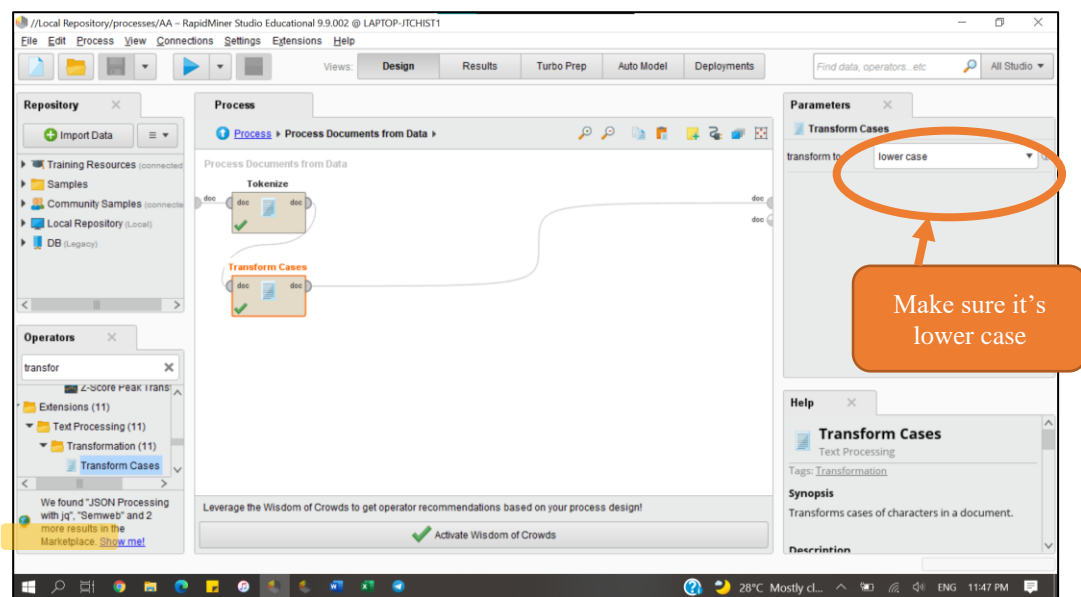
**Input** :



Figure 10 : Add 'Transform Cases' operator

**Output** :



Figure 11  : Result of transform cases process

## Step 6 : Filter Stopwords

After transforming the dataset into lower cases, next step is filtering all stopwords. Stopwords are common words like "a", "the", "is", "are" and etc. We removed those common words so that we can focus on important words. In this, I used the filter stopwords by dictionary and manually add a text file that contains a thousand of stopwords that has been listed. To do so, I needed to add the 'Open File' operator. After the execution, a list of words without stopwords has been generated and the total of words generated are 29 after going through the Filtering Stopwords process.
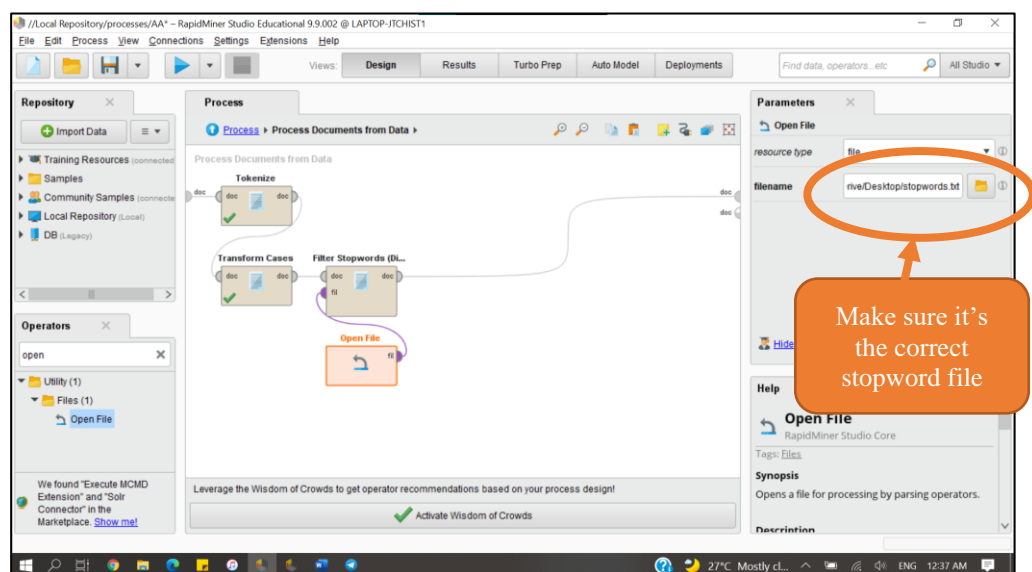
**Input :**



Figure 12 : Add the 'Filter Stopword by dictionary' operator and the stopword file manually
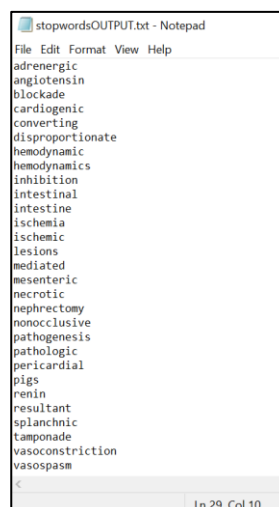
**Output :**



Figure 12 : List of words without stopwords

## Step 7 : Generate n grams (terms)

Next, we need to generate n-Grams. Clearly, language has a successive nature, thus the request in which words show up in the content matters a great deal. This operator permits us to comprehend the context of a sentence regardless of whether there are a few words missing. To do this, drag the 'Generate n-Grams' operator onto the process.
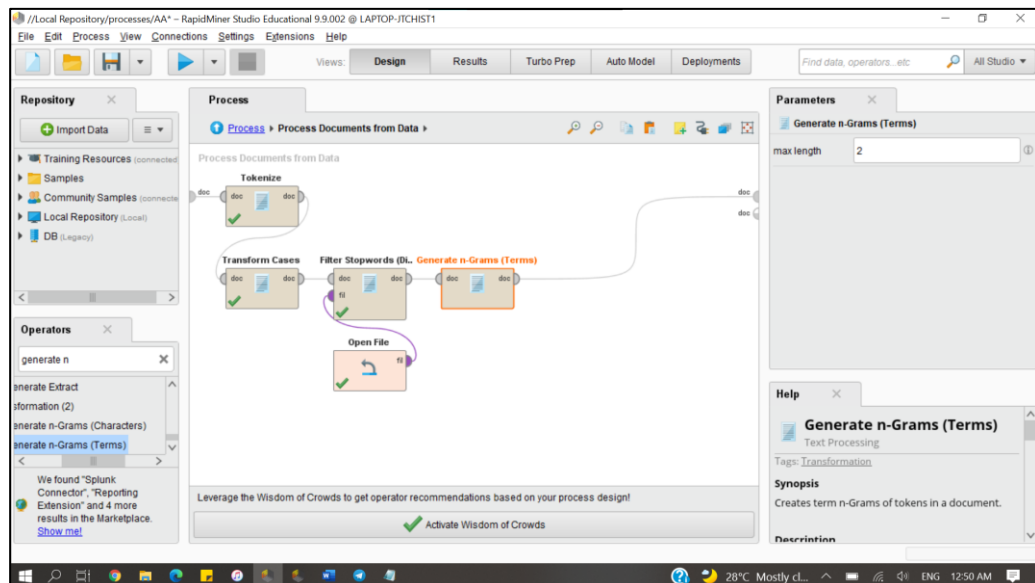
**Input :**



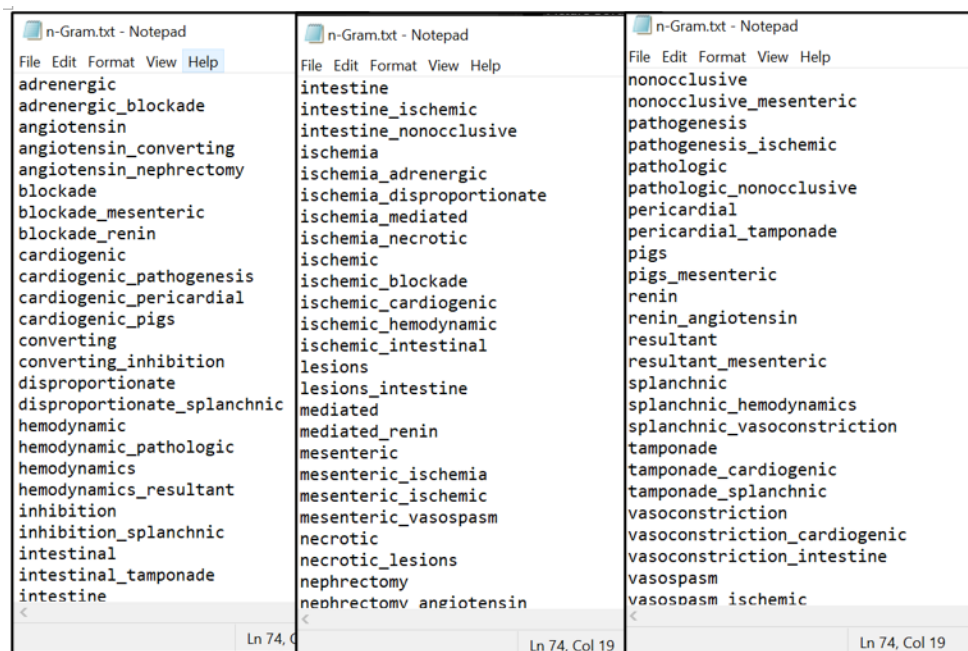Figure 13 : Add the 'Generate n-Gram' operator

**Output** :



Figure 14 : The result after going through the Generate n-Gram process

## Step 8 : Filter Tokens by POS Tags

The next step is Filter Tokens by POS Tags. A POS tag (or part-of-speech tag) is a special label assigned to each token (word) in a text corpus to indicate the part of speech and often also other grammatical categories such as tense, number (plural/singular), case etc. To complete this, drag the 'Filter Tokens by POS Tags' operator and set the expression to N.*
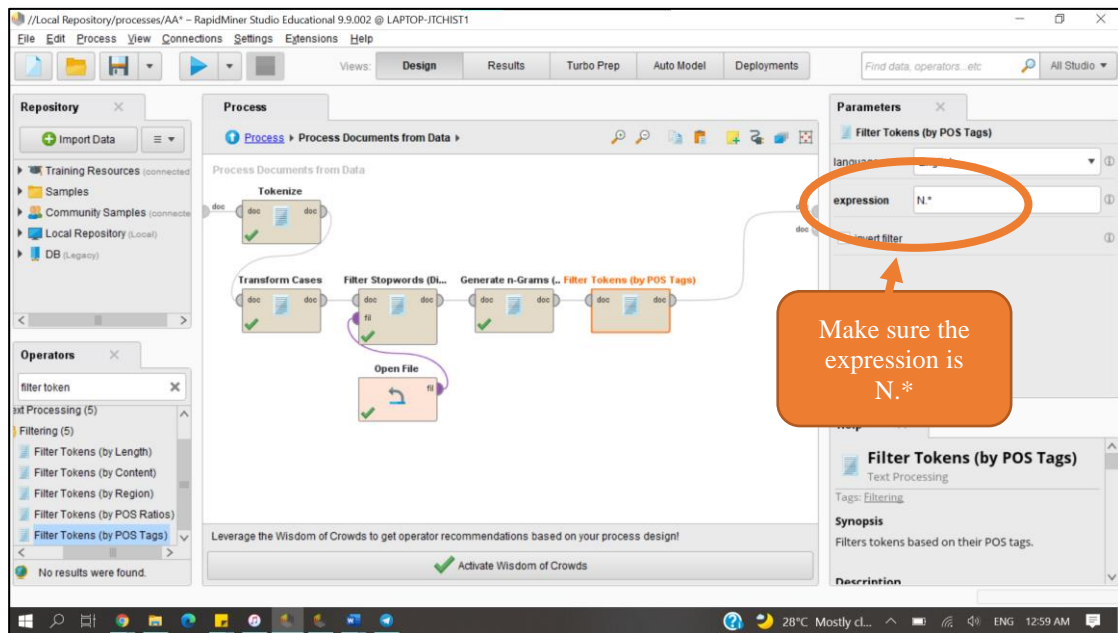
**Input** :



Figure 15 : Add the 'Filter Tokens by POS Tags' operator

**Output** :



Figure 16 : The list of word generated after filter token process has been executed

11

## Step 9 : Stem (Porter)

Then, I needed to do the stemming step. The Porter stemming algorithm (or 'Porter stemmer') is a process for removing the commoner morphological and in flexional endings from words in English. To do this, simply drag the 'Stem(Porter)' operator.

**Input** :



Figure 17 : Add the 'Stem (Porter)' operator into the process

**Output** :



Figure 18 : A list of words generated after the execution of Stem (Porter) process

## Step 10 : Wordlist to Data

## Input :



**Figure 19 : Add the 'Wordlist to Data' operator**

## Output :



**Figure 20 : The final result**
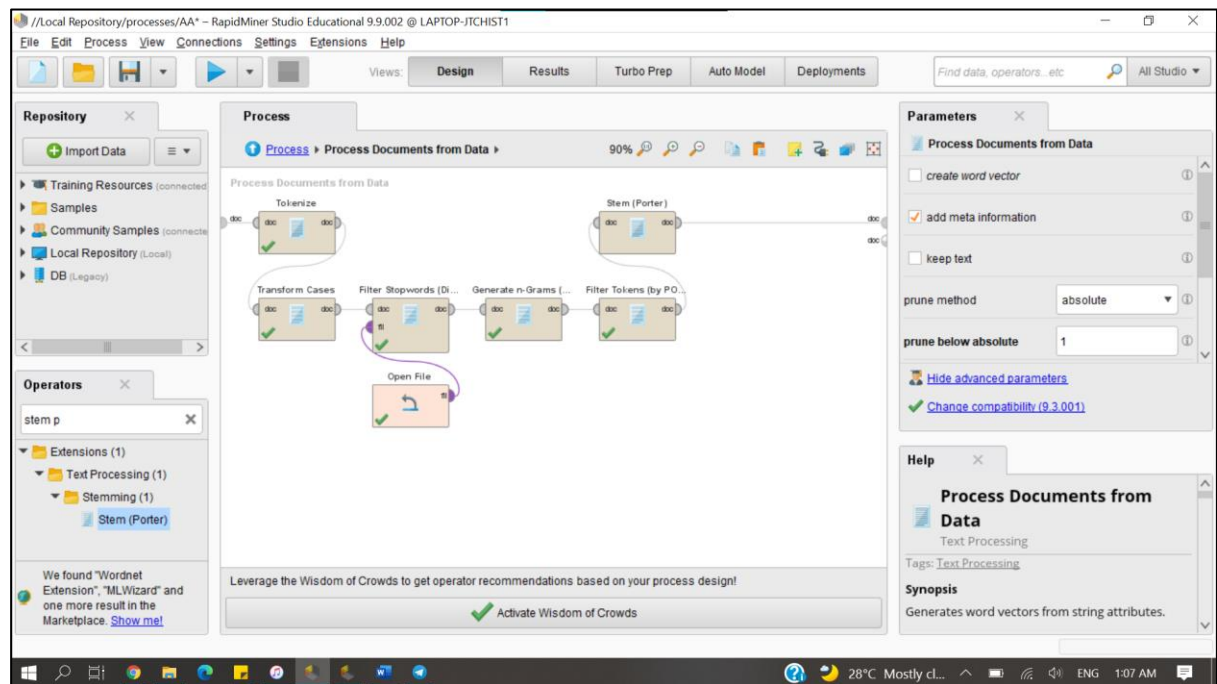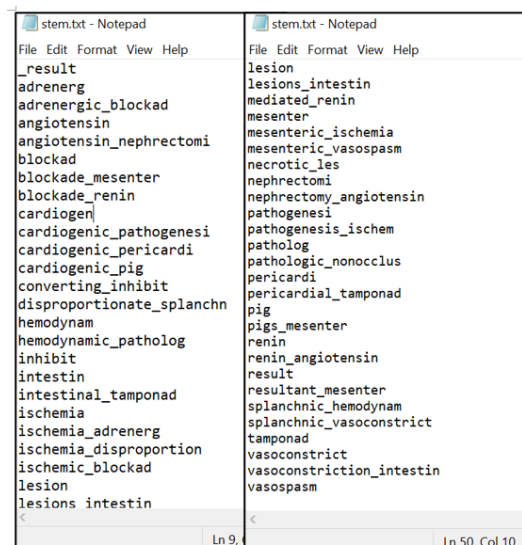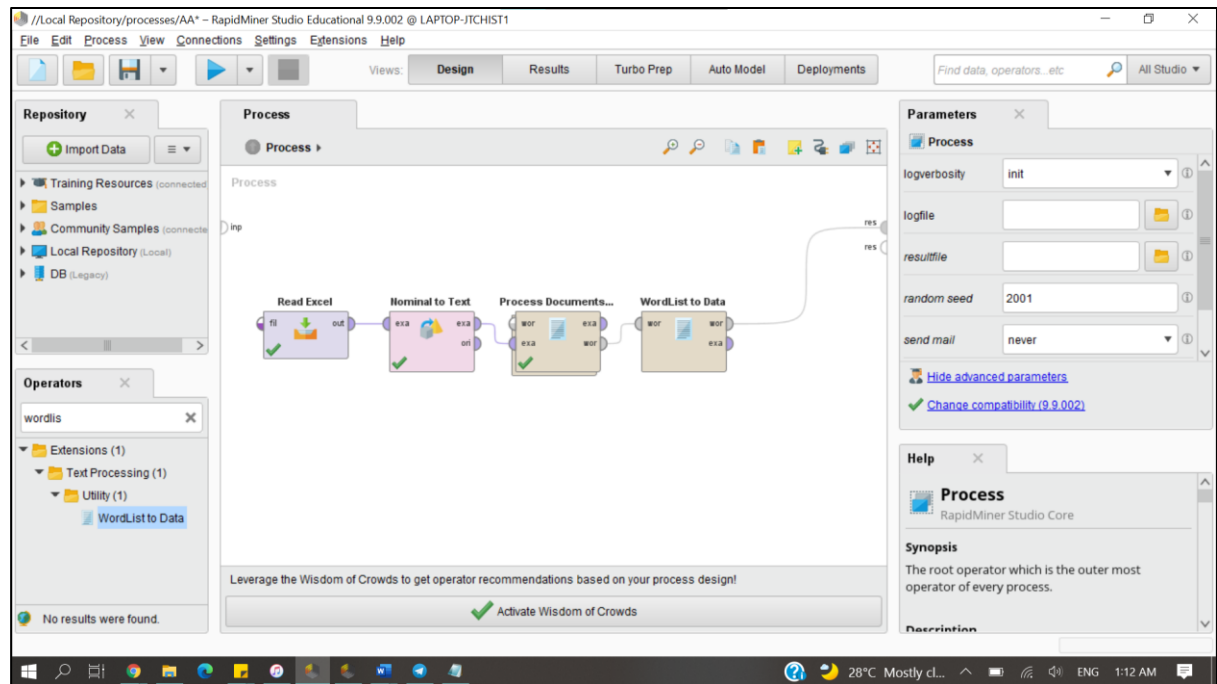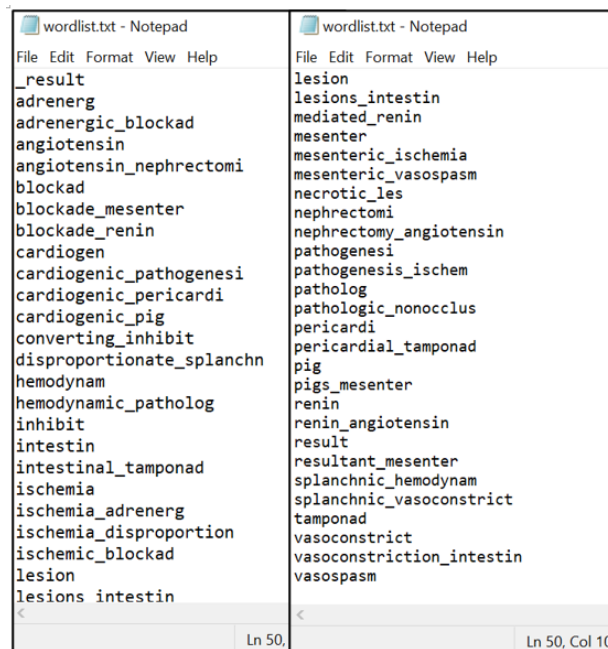
## 1.2 Data Mining Task

### 1.2.1 Supervised Learning

**a)**

The one supervised learning algorithm that I can apply for appendix 2 is Classification. There are 2 common algorithms for the supervised which are Regression and Classification. Numerical data is made of numbers and it has two more categories that falls under it which are continuous and discrete. To use the Regression, we need the data in continuous target variable. Figure below shows the example of continuous data.

| x0 | x1 | x2 | x3 | x4 | x5 |
|---|---|---|---|---|---|
| 1 | 1.06 | 9.2 | 151 | 54.4 | 1.6 |
| 2 | 0.89 | 10.3 | 202 | 57.9 | 2.2 |
| 3 | 1.43 | 15.4 | 113 | 53 | 3.4 |
| 4 | 1.02 | 11.2 | 168 | 56 | 0.3 |
| 5 | 1.49 | 8.8 | 192 | 51.2 | 1 |
| 6 | 1.32 | 13.5 | 111 | 60 | -2.2 |
| 7 | 1.22 | 12.2 | 175 | 67.6 | 2.2 |
| 8 | 1.1 | 9.2 | 245 | 57 | 3.3 |
| 9 | 1.34 | 13 | 168 | 60.4 | 7.2 |
| 10 | 1.12 | 12.4 | 197 | 53 | 2.7 |
| 11 | 0.75 | 7.5 | 173 | 51.5 | 6.5 |
| 12 | 1.13 | 10.9 | 178 | 62 | 3.7 |

Figure 21 : Continuous data

While on the other hand, we also have Categorical data. Categorical data is made of words. Under it. It has two more categories which are ordinal and nominal. Classification algorithm is used to interpret the nominal data. Just like the appendix 2, it its made of words. Therefore, the most suitable supervised learning algorithm is Classification. Figure below shows an example of nominal data.

| | A | B | C |
|---|---|---|---|
| 1 | Order ID | Product Name | Feedback |
| 2 | Order 1 | Sugar | Positive |
| 3 | Order 1 | Bread | Positive |
| 4 | Order 2 | Sugar | Negative |
| 5 | Order 2 | Bread | Positive |
| 6 | Order 2 | Rice | Negative |
| 7 | Order 2 | Soda | Positive |
| 8 | Order 3 | Sugar | Negative |
| 9 | Order 3 | Bread | Negative |
| 10 | Order 3 | Rice | Negative |
| 11 | Order 3 | Soda | Positive |

Figure 22 : Nominal data

**b)**

## Step 1 : Read Excel

To prepare the dataset for training, I read my dataset 29 that I already converted into an excel file by using the Read Excel operator.
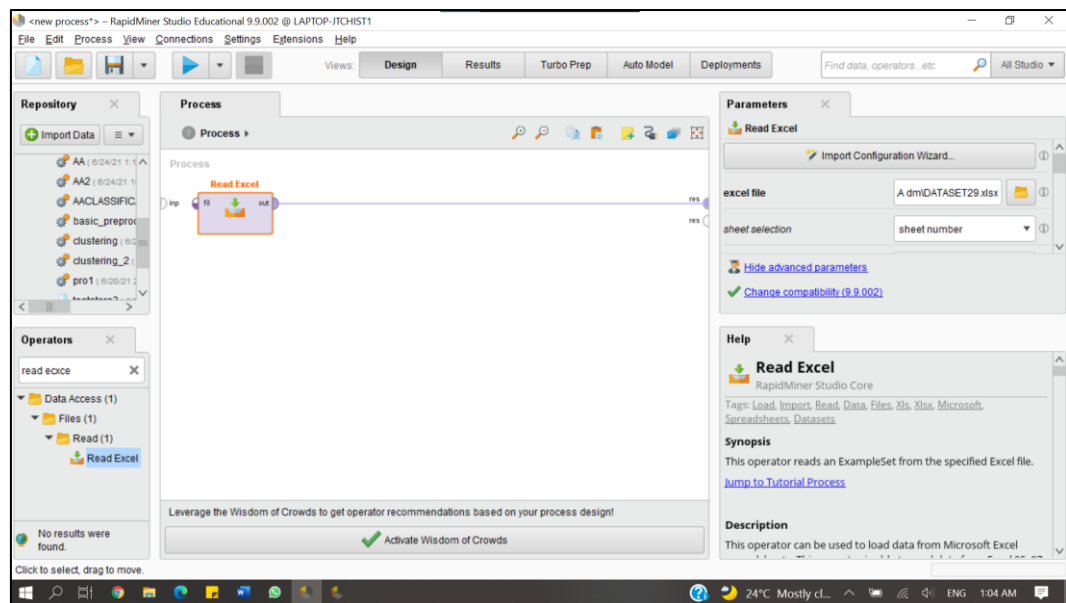
## Input :



Figure 23 : Read the data from excel

**Output :**

| Row No. | FILENAME | disease | ABSTRACT |
|---|---|---|---|
| 1 | 13918 | BRADYCARDIA | Surgical treatment of pediatric cardiac arrhyth... |
| 2 | 14824 | BRADYCARDIA | Comparative survival following permanent ven... |
| 3 | 16526 | BRADYCARDIA | Complete sinoatrial block in two patients with ... |
| 4 | 275357 | CORONARY DISEASE | Enhanced utilization of exogenous glucose im... |
| 5 | 275546 | CORONARY DISEASE | Myocardial amiodarone and desethylamiodar... |
| 6 | 275547 | CORONARY DISEASE | Effects of benazepril and metoprolol OROS al... |
| 7 | 18126 | HEART ANEURYSM | Atrial septal aneurysms in infants and children. |
| 8 | 27166 | HEART ANEURYSM | Submitral left ventricular aneurysms. Correctio... |
| 9 | 9419 | MYOCARDIAL DISEASES | Altered norepinephrine turnover and metaboli... |
| 10 | 10031 | MYOCARDIAL DISEASES | Pathophysiology and pathogenesis of stunne... |

Figure 24 : Data that has been read using read excel

## Step 2 : Nominal To Text

Next, I added the Nominal to Text operator. The Nominal to Text operator converts all nominal attributes to string attributes.
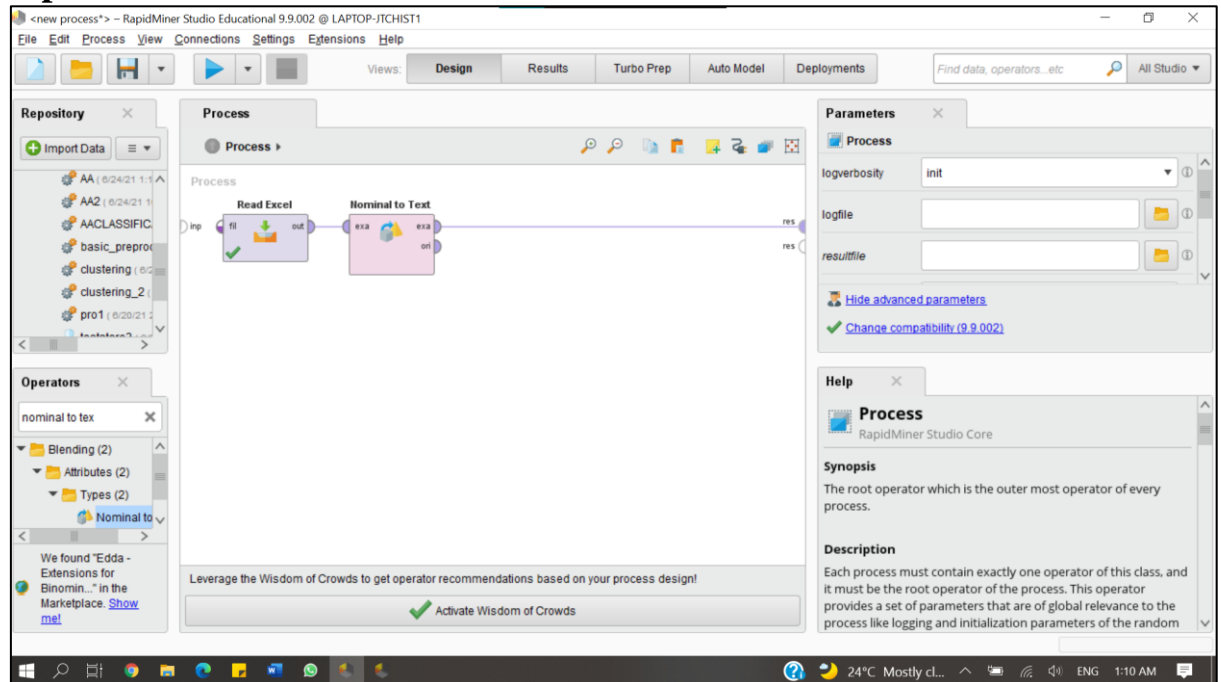
## Input :



Figure 25 : Add the Nominal To Text operator

## Output :

| Row No. | FILENAME | disease | ABSTRACT |
|---------|----------|---------|----------|
| 1 | 13918 | BRADYCARDIA | Surgical treatment of pediatric cardiac arrhyth... |
| 2 | 14824 | BRADYCARDIA | Comparative survival following permanent ven... |
| 3 | 16526 | BRADYCARDIA | Complete sinoatrial block in two patients with ... |
| 4 | 275357 | CORONARY DISEASE | Enhanced utilization of exogenous glucose im... |
| 5 | 275546 | CORONARY DISEASE | Myocardial amiodarone and desethylamiodar... |
| 6 | 275547 | CORONARY DISEASE | Effects of benazepril and metoprolol OROS al... |
| 7 | 18126 | HEART ANEURYSM | Atrial septal aneurysms in infants and children. |
| 8 | 27166 | HEART ANEURYSM | Submitral left ventricular aneurysms. Correctio... |
| 9 | 9419 | MYOCARDIAL DISEASES | Altered norepinephrine turnover and metaboli... |
| 10 | 10031 | MYOCARDIAL DISEASES | Pathophysiology and pathogenesis of stunne... |

Figure 26 : The data after has changed to text

## Step 3 : Select Attribute

Next, I will select the attributes. From here I selected ABSTRACT and disease as the chosen attributes.
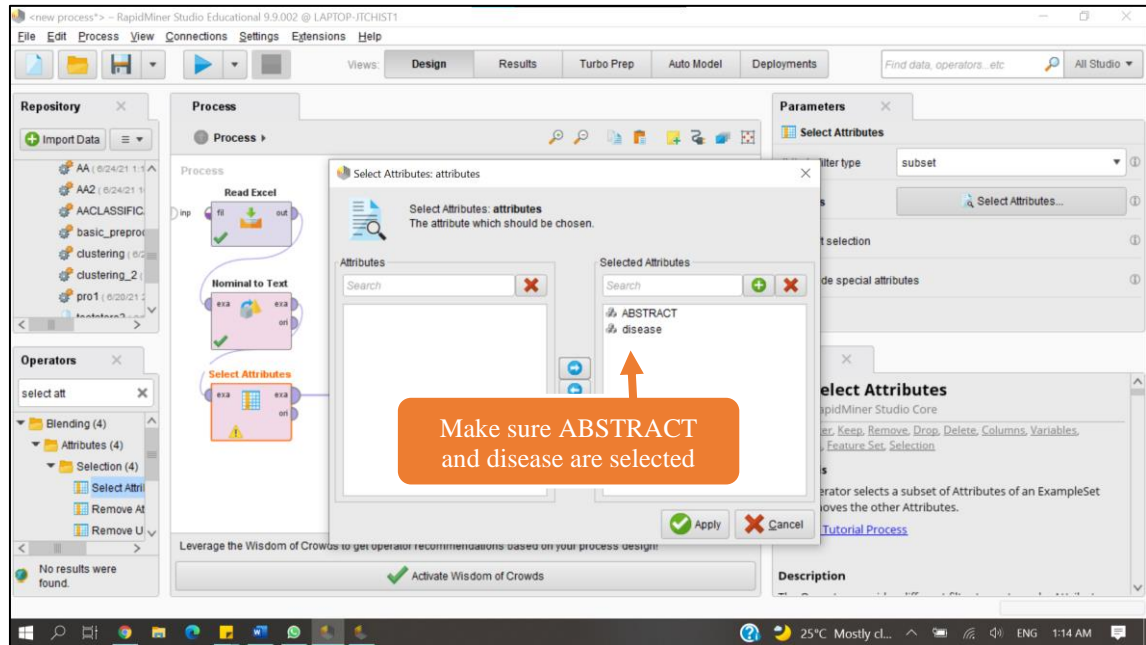
## Input :



Figure 27 : Select Attribute

## Output :



| Row No. | disease | ABSTRACT |
|---|---|---|
| 1 | BRADYCARDIA | Surgical treatment of pediatric cardiac arrhythmia. |
| 2 | BRADYCARDIA | Comparative survival following permanent ventricul... |
| 3 | BRADYCARDIA | Complete sinoatrial block in two patients with brad... |
| 4 | CORONARY DISEASE | Enhanced utilization of exogenous glucose improv... |
| 5 | CORONARY DISEASE | Myocardial amiodarone and desethylamiodarone c... |
| 6 | CORONARY DISEASE | Effects of benazepril and metoprolol OROS alone ... |
| 7 | HEART ANEURYSM | Atrial septal aneurysms in infants and children. |
| 8 | HEART ANEURYSM | Submitral left ventricular aneurysms. Correction by ... |
| 9 | MYOCARDIAL DISEASES | Altered norepinephrine turnover and metabolism i... |
| 10 | MYOCARDIAL DISEASES | Pathophysiology and pathogenesis of stunned my... |

Figure 28 : The result after choosing specific attributes

## Step 4 : Set Role

Then, I set the role for disease as label. Label is a special role. An Attribute with the id role acts as an identifier for the Examples. It should be unique for all Examples.
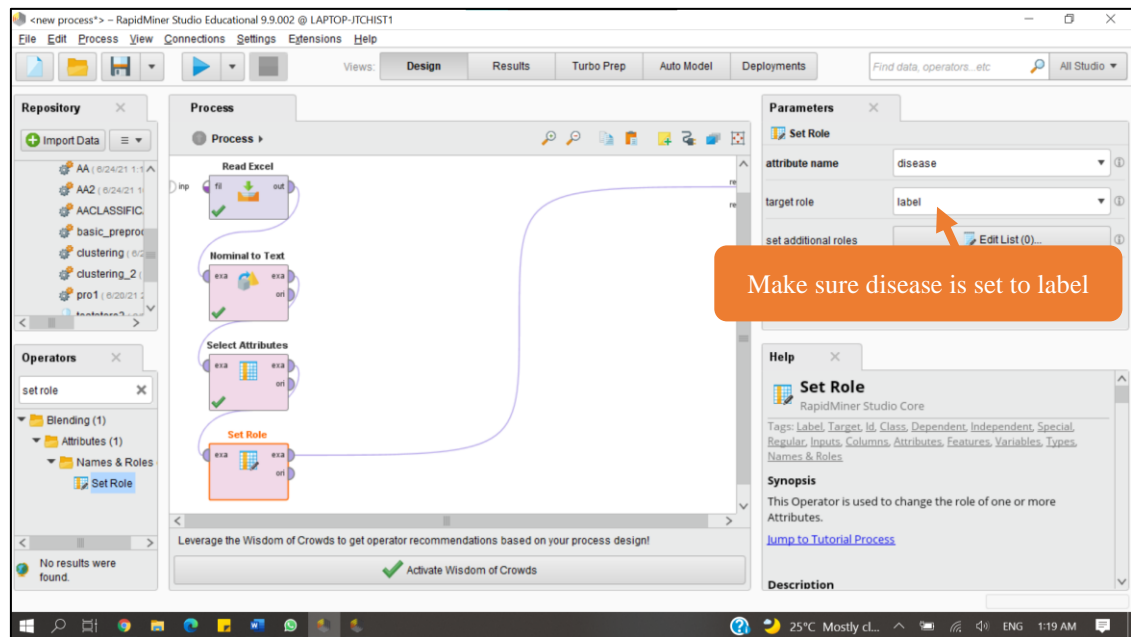
## Input :



Figure 29 : Set the role label for Disease

## Output :



Figure 30 : The disease after has been set role as label

## Step 5 : Split Data

Continue, I split data into training and testing dataset by 0.8 (80%) and 0.2 (20%) respectively. Splitting data into training and testing sets is a significant piece of assessing data mining models. By utilizing comparable information for training and testing, we can limit the impacts of information errors and better comprehend the qualities of the model. [2] After the completion I keep the testing data into a file data named teststore1.
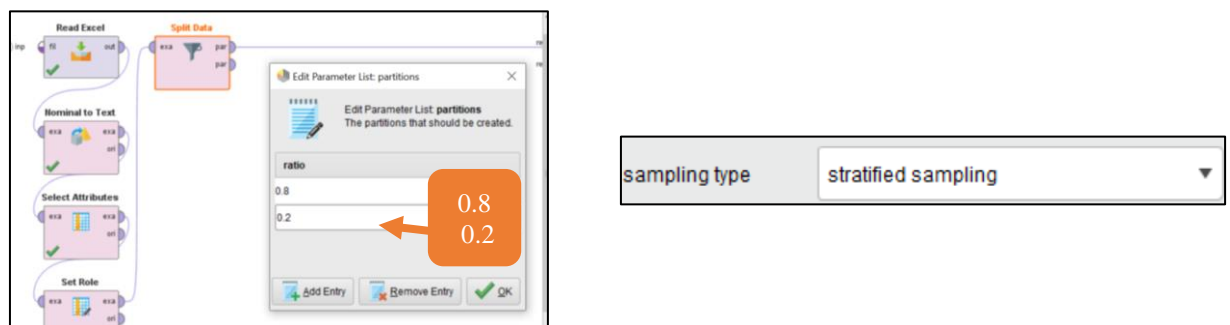
## Input :



Figure 31 : Split Data into 0.8 and 0.2 and change the sampling type

## Output :

| Row No. | disease | ABSTRACT |
|---------|---------|----------|
| 1 | BRADYCARDIA | Comparative survival following permanent ventricular and dual-cham... |
| 2 | BRADYCARDIA | Complete sinoatrial block in two patients with bradycardia-tachycardi... |
| 3 | CORONARY DISEASE | Enhanced utilization of exogenous glucose improves cardiac function... |
| 4 | CORONARY DISEASE | Effects of benazepril and metoprolol OROS alone and in combination ... |
| 5 | HEART ANEURYSM | Atrial septal aneurysms in infants and children. |
| 6 | HEART ANEURYSM | Submitral left ventricular aneurysms. Correction by a new transatrial a... |
| 7 | MYOCARDIAL DISEASES | Altered norepinephrine turnover and metabolism in diabetic cardiomy... |
| 8 | MYOCARDIAL DISEASES | Pathophysiology and pathogenesis of stunned myocardium. Depress... |

Figure 32 : The training dataset (80%)

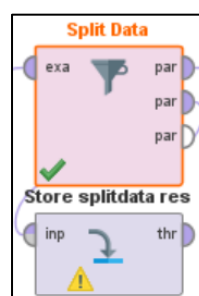| Row No. | disease | ABSTRACT |
|---------|---------|----------|
| 1 | BRADYCARDIA | Surgical treatment of pediatric car... |
| 2 | CORONARY DISEASE | Myocardial amiodarone and deset... |

Figure 33 : The testing dataset (20%)



Figure 34 : Store the result into a data file

## Step 6 : Preprocessing

Before starting any process, I should do data cleaning to my dataset. Therefore, I added Process Document from Data operator and do all preprocessing steps inside the operator.
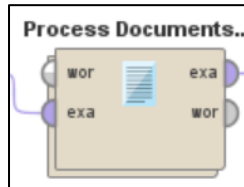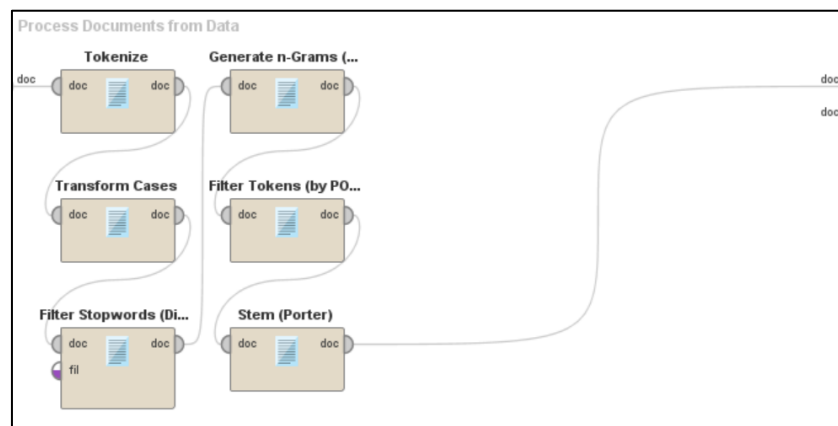
## Input :



Figure 35 : Add Process Documents



Figure 36 : Perform Preprocessing steps

## Output :



| Word | .. | Total O... | Document Occurences | BRADYCARDIA | CORONARY DISEASE | HEART ANEURYSM | MYOCARDIAL DISEASES |
|------|----|-----------|---------------------|-------------|------------------|----------------|---------------------|
| aneurysm | . | 12 | 2 | 0 | 0 | 12 | 0 |
| atrial | . | 11 | 3 | 6 | 0 | 5 | 0 |
| atrium | . | 2 | 2 | 1 | 0 | 1 | 0 |
| bradycar... | . | 7 | 2 | 7 | 0 | 0 | 0 |
| bradycar... | . | 4 | 2 | 4 | 0 | 0 | 0 |
| dysfunct | . | 7 | 3 | 3 | 2 | 0 | 2 |
| ischemia | . | 2 | 2 | 0 | 1 | 0 | 1 |
| markedli | . | 2 | 2 | 0 | 1 | 0 | 1 |
| maxim | . | 6 | 2 | 0 | 3 | 0 | 3 |
| sinoatri | . | 3 | 2 | 3 | 0 | 0 | 0 |
| sinoatria... | . | 2 | 2 | 2 | 0 | 0 | 0 |
| sinu | . | 10 | 2 | 10 | 0 | 0 | 0 |
| sinus_br... | . | 3 | 2 | 3 | 0 | 0 | 0 |
| tachycar... | . | 5 | 2 | 5 | 0 | 0 | 0 |

Figure 37 : Data after going through preprocessing

## Step 7 : Cross Validation

Next, I added Cross Validation operator. Cross-validation is fundamentally utilized in applied machine learning to assess the ability of a machine learning model on unseen data. That is, to utilize a restricted example to assess how the model is relied upon to act overall when used to make expectations on data not utilized during the preparation of the model. [3]
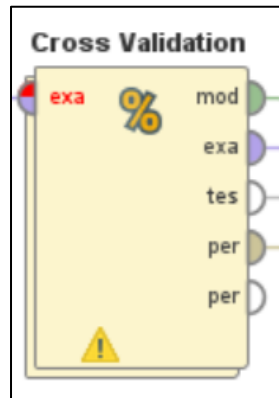
**Input :**


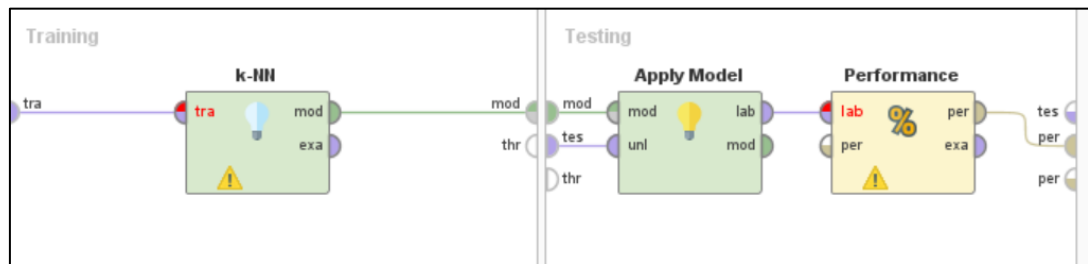
Figure 38 : Add Cross Validation operator



Figure 39 : Inside Cross Validation

**Output :**

| accuracy: 25.00% | | | | | |
|---|---|---|---|---|---|
| | true BRADYCARDIA | true CORONARY D... | true HEART ANEU... | true MYOCARDIAL ... | class precision |
| pred. BRADYCARD... | 1 | 0 | 1 | 0 | 50.00% |
| pred. CORONARY ... | 0 | 0 | 0 | 2 | 0.00% |
| pred. HEART ANE... | 1 | 0 | 1 | 0 | 50.00% |
| pred. MYOCARDIA... | 0 | 2 | 0 | 0 | 0.00% |
| class recall | 50.00% | 0.00% | 50.00% | 0.00% | |

Figure 40 : The accuracy

```
PerformanceVector

PerformanceVector:
accuracy: 0.00%
ConfusionMatrix:
True:    BRADYCARDIA      CORONARY DISEASE      HEART ANEURYSM  MYOCARDIAL DISEASES
BRADYCARDIA :    0        0        2        0
CORONARY DISEASE :       0        0        0        2
HEART ANEURYSM :         2        1        0        0
MYOCARDIAL DISEASES:     0        1        0        0
kappa: -0.333
ConfusionMatrix:
True:    BRADYCARDIA      CORONARY DISEASE      HEART ANEURYSM  MYOCARDIAL DISEASES
BRADYCARDIA :    0        0        2        0
CORONARY DISEASE :       0        0        0        2
HEART ANEURYSM :         2        1        0        0
MYOCARDIAL DISEASES:     0        1        0        0
```

Figure 41 : Performance Vector

## Testing

## Step 1 : Get the data and test

To start the testing, I fetched my data and testing dataset from the file that I have stored before during the process Training.
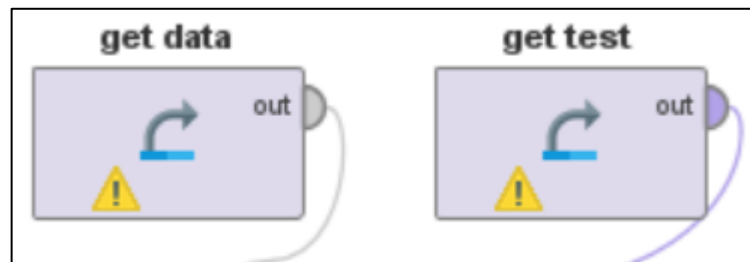
## Input :



Figure 42 : Get Data and Test

## Output :



Figure 43 : Result for get Test

| Word | Attribut... | Total O... | Docum... | BRADY... | CORON... | HEART ... | MYOCA... |
|------|-------------|-----------|----------|----------|----------|-----------|----------|
| aneurysm | aneurysm | 12 | 2 | 0 | 0 | 12 | 0 |
| atrial | atrial | 11 | 3 | 6 | 0 | 5 | 0 |
| atrium | atrium | 2 | 2 | 1 | 0 | 1 | 0 |
| bradycardia | bradycar... | 7 | 2 | 7 | 0 | 0 | 0 |
| bradycardia_tachycardia | bradycar... | 4 | 2 | 4 | 0 | 0 | 0 |
| dysfunct | dysfunct | 7 | 3 | 3 | 2 | 0 | 2 |
| ischemia | ischemia | 2 | 2 | 0 | 1 | 0 | 1 |
| markedli | markedli | 2 | 2 | 0 | 1 | 0 | 1 |
| maxim | maxim | 6 | 2 | 0 | 3 | 0 | 3 |
| sinoatri | sinoatri | 3 | 2 | 3 | 0 | 0 | 0 |
| sinoatrial_bradycardia | sinoatria... | 2 | 2 | 2 | 0 | 0 | 0 |

Figure 44 : Result for get Data

## Step 2 : Preprocess

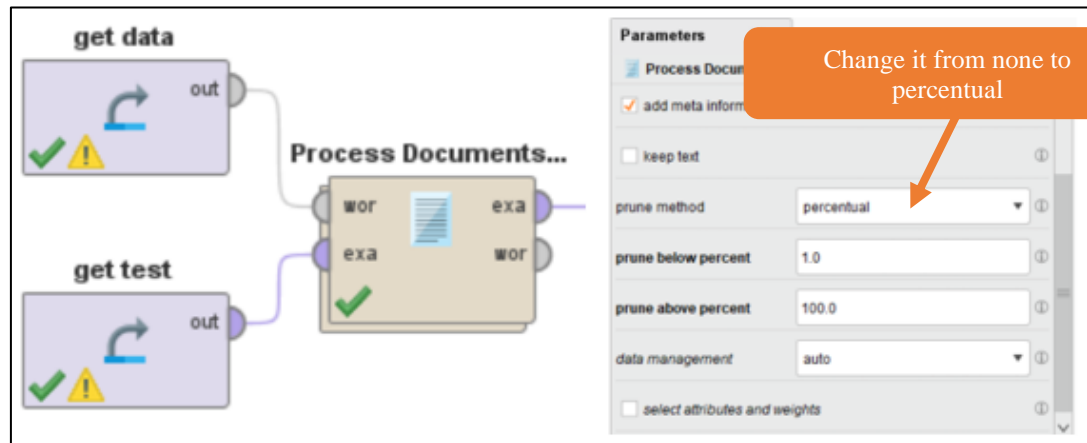To double the secure, I preprocessed both data again and the prune method is percentual.
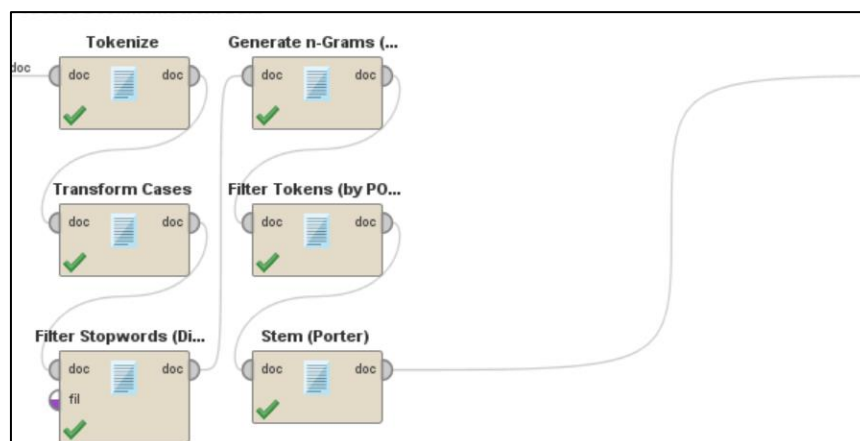
## Input :


Figure 45 : Add Process Document


Figure 46 : Inside of the process document

## Output :

| Row No. | disease | aneurysm | atrial | atrium | bradycardia | bradycardia... | dysfunct | ischemia | marke |
|---------|---------|----------|--------|--------|-------------|----------------|----------|----------|-------|
| 1 | BRADYCARD... | 0 | 0 | 0 | 0.707 | 0 | 0 | 0 | 0 |
| 2 | CORONARY ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Figure 47 : The result after going through preprocess

## Step 3 : Set Role

Next, I set role for disease as the label. This step is important. Setting an attribute to label is a very important role. An Attribute with the id role goes about as an identifier for the Examples. It should be unique for all Examples.
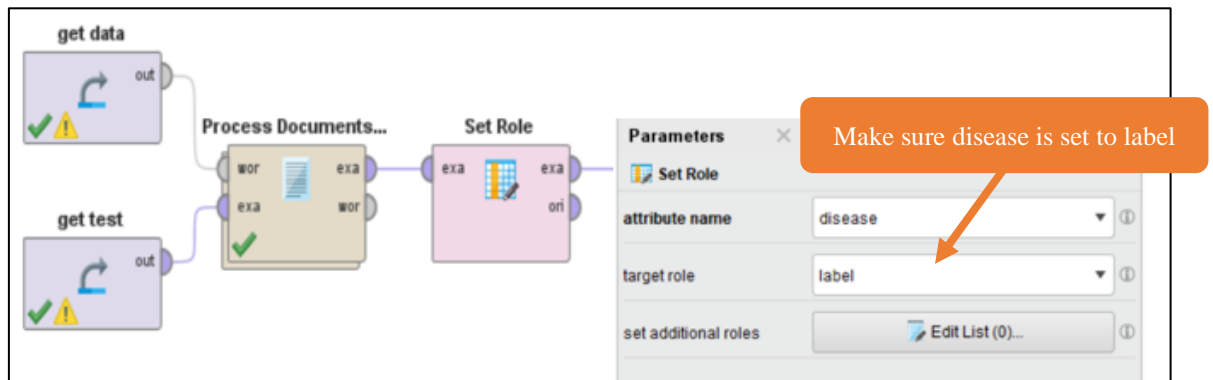
## Input :



Figure 48 : Set Role for attribute disease

## Output :

| Row No. | disease | aneurysm | atrial | atrium | bradycardia | bradycardia... | dysfunct | ischemia | marke |
|---------|---------|----------|--------|--------|-------------|----------------|----------|----------|-------|
| 1 | BRADYCARD... | 0 | 0 | 0 | 0.707 | 0 | 0 | 0 | 0 |
| 2 | CORONARY ... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Figure 49 : The result after set role

## Step 4 : Apply Model

Lastly, I will apply the model. I added the Apply Model operator. As I add the operator, it will ask for the model. Therefore, I imported the model that has been generated during the Training process that I kept in a file by using the Retrieve operator.
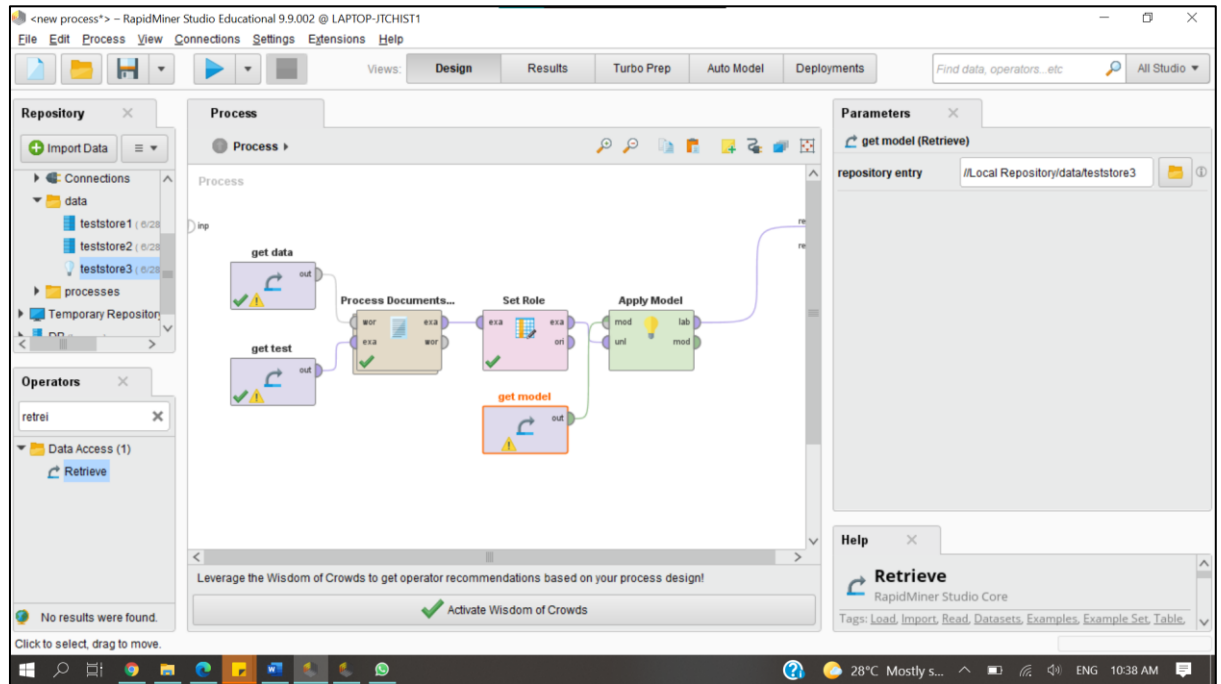
## Input :



Figure 50 : Add Apply Model operator and get Model operator

## Output :

| Row No. | disease | prediction(disease) | confidence(... | confidence(... | confidence(... | confidence(... | aneurysm | atri |
|---------|---------|---------------------|----------------|----------------|----------------|----------------|----------|------|
| 1 | BRADYCARDIA | BRADYCARDIA | 0.704 | 0.296 | 0 | 0 | 0 | 0 |
| 2 | CORONARY DISEASE | HEART ANEURYSM | 0.354 | 0 | 0.646 | 0 | 0 | 0 |

Figure 51 : The final result of the classification



Figure 52 : KNNClassification

## c) Performance Measures

| Row No. | disease | prediction(disease) | confidence(... | confidence(... | confidence(... | confidence(... | aneurysm | atri |
|---------|---------|---------------------|----------------|----------------|----------------|----------------|----------|------|
| 1 | BRADYCARDIA | BRADYCARDIA | 0.704 | 0.296 | 0 | 0 | 0 | 0 |
| 2 | CORONARY DISEASE | HEART ANEURYSM | 0.354 | 0 | 0.646 | 0 | 0 | 0 |

Figure 53 : The result for Testing dataset using KNN algorithm

accuracy: 25.00%

| | true BRADYCARDIA | true CORONARY D... | true HEART ANEU... | true MYOCARDIAL ... | class precision |
|---|---|---|---|---|---|
| pred. BRADYCARD... | 1 | 0 | 1 | 0 | 50.00% |
| pred. CORONARY ... | 0 | 0 | 0 | 2 | 0.00% |
| pred. HEART ANE... | 1 | 0 | 1 | 0 | 50.00% |
| pred. MYOCARDIA... | 0 | 2 | 0 | 0 | 0.00% |
| class recall | 50.00% | 0.00% | 50.00% | 0.00% | |

Figure 54 : The result for Training dataset using KNN algorithm

Based on figure 54, the total item in the table could total up to 8, which is equal to the training dataset. It has been shown that the accuracy of the training dataset is 25%. To understand the training result better, we should look at BRADYCARDIA. It predicted BRADYCARDIA to be true is only 1 while the another one it predicted it to be HEART ANEURYSM. This could cause the precision to be 50% while recall is also 50%. Even though the accuracy is low, we can't say that the whole classification process using KNN is totally wrong. The accuracy is actually depends on number of data used, the less data been used will cause the result to be a little inaccurate and insufficient.

## 1.2.2 Unsupervised Learning Algorithm

Unsupervised learning is likewise an exceptionally normal kind of machine learning. It contrasts from regulated learning in that the data that has no label (unlabeled). One unsupervised learning algorithm that I can apply for given appendix 2 is Clustering. There are many algorithms that falls under the unsupervised learning algorithm such as clustering, association and dimensionality reduction. Clustering is gathering a set of items in such a way that objects in a same cluster are more comparable than to those objects having a place with different cluster. While, association rules are tied in with discovering relationship among things inside huge business data sets. Clustering algorithms are for the most part utilized when we need to make the clusters dependent on the characteristic of the data focuses. For both algorithms, the target variable is not available as the data that is going to use is the unlabeled.



Figure 55: Visualization of data [4]

**e)**

## Step 1 : Read Excel

The very first step is to read the dataset. For this question, I needed to refer appendix 2 and my dataset is dataset 29. Therefore, I changed the format from .docx to .xlsx so that the tool can read it easily by using the 'Read Excel' operator.

**Input** :



Figure 56 : Read excel operator

**Output** :

| Row No. | disease | ABSTRACT |
|---------|---------|----------|
| 1 | BRADYCARDIA | Surgical treatment of pediatric cardiac arrhythmia. |
| 2 | BRADYCARDIA | Comparative survival following permanent ventricular and dual-chamber p... |
| 3 | BRADYCARDIA | Complete sinoatrial block in two patients with bradycardia-tachycardia syn... |
| 4 | CORONARY DISEASE | Enhanced utilization of exogenous glucose improves cardiac function in hy... |
| 5 | CORONARY DISEASE | Myocardial amiodarone and desethylamiodarone concentrations in patient... |
| 6 | CORONARY DISEASE | Effects of benazepril and metoprolol OROS alone and in combination on ... |
| 7 | HEART ANEURYSM | Atrial septal aneurysms in infants and children. |
| 8 | HEART ANEURYSM | Submitral left ventricular aneurysms. Correction by a new transatrial appro... |
| 9 | MYOCARDIAL DISEASES | Altered norepinephrine turnover and metabolism in diabetic cardiomyopat... |
| 10 | MYOCARDIAL DISEASES | Pathophysiology and pathogenesis of stunned myocardium. Depressed C... |

Figure 57 : The data that has been read by the tool

## Step 2 : Nominal to Text

Next, I added the Nominal to Text operator. The Nominal to Text operator converts all nominal attributes to string attributes.

**Input** :



Figure 58 : Nominal to Text operator

**Output** :



| Row No. | disease | ABSTRACT |
|---------|---------|----------|
| 1 | BRADYCARDIA | Surgical treatment of pediatric cardiac arrhythmia. |
| 2 | BRADYCARDIA | Comparative survival following permanent ventricular and d... |
| 3 | BRADYCARDIA | Complete sinoatrial block in two patients with bradycardia-t... |
| 4 | CORONARY DISEASE | Enhanced utilization of exogenous glucose improves cardi... |
| 5 | CORONARY DISEASE | Myocardial amiodarone and desethylamiodarone concentr... |
| 6 | CORONARY DISEASE | Effects of benazepril and metoprolol OROS alone and in co... |
| 7 | HEART ANEURYSM | Atrial septal aneurysms in infants and children. |
| 8 | HEART ANEURYSM | Submitral left ventricular aneurysms. Correction by a new tr... |
| 9 | MYOCARDIAL DISEASES | Altered norepinephrine turnover and metabolism in diabeti... |
| 10 | MYOCARDIAL DISEASES | Pathophysiology and pathogenesis of stunned myocardiu... |

Figure 59 : The data has been changed from nominal to text

## Step 3 : Select attribute

In this step, I needed to select the needed atrtribute only. Therefore I chose ABSTRACT as I am about to process the long text data.

## Input :



Figure 60 : Choosing the Abstract as attribute

## Output :

| Row No. | ABSTRACT |
|---------|----------|
| 1 | Surgical treatment of pediatric cardiac arrhthmia. |
| 2 | Comparative survival following permanent ventricular an... |
| 3 | Complete sinoatrial block in two patients with bradycardi... |
| 4 | Enhanced utilization of exogenous glucose improves car... |
| 5 | Myocardial amiodarone and desethylamiodarone conce... |
| 6 | Effects of benazepril and metoprolol OROS alone and in ... |
| 7 | Atrial septal aneurysms in infants and children. |
| 8 | Submitral left ventricular aneurysms. Correction by a ne... |
| 9 | Altered norepinephrine turnover and metabolism in diab... |
| 10 | Pathophysiology and pathogenesis of stunned myocardi... |

Figure 61 : The latest data after choosing attribute abstract

## Step 4 : Preprocessing

To do preprocessing, I needed to add the Process Documents from Data operator first. And then proceed the basic preprocessing in the Process Documents. Almost all process, preprocessing is the most basic and needed to be done first before we get deeper into other processes.

## Input :



Figure 62 : Add Process Documents operator



Figure 63 : All preprocessing basic step

## Output :



Figure 64 : The list of word generated after going through the preprocessing

## Step 5 : K Means

To apply the clustering, I added the K-Means Clustering operator. This Operator performs clustering utilizing the k-means algorithm. Clustering assembles examples which are almost likely the same with one another.

## Input :



Figure 65 : Add the K-Means operator

## Output :



Figure 66 : Example Set



Figure 67 : Cluster Model

## Step 6 : Cluster Model Visualizer

To visualize and understand the cluster model better, I added Cluster Model Visualizer. There will be no parameter needed.

### Input :



Figure 68 : Add Cluster Model Visualizer

### Output :



Figure 69 : Output Cluster Model Visualizer

## f) Performance Measures



Figure 70 : Output Cluster Model Visualizer

Based on the figure above, it shows the final result for the Clustering process and it is visualized using the Cluster Model Operator. The total number of clusters is 5 (cluster 0 – cluster 4) since I set the value of K = 5 during the process of K-Means. To understand this better, from the figure, in cluster 3, aneurysm, atrial and atrium are been clustered together since they own almost the same characteristic. But we can see that the word aneurysm is 400% larger than atrial (189.24%) and atrium (150.35%). This simply means that aneurysm is the most popular word in the cluster 3.

# 2.0 PART 2

(a) Throughout completing this alternative assessment, I found that I learnt a lot of new things specifically in Data Mining. Before doing this alternative assessment, I thought it's very hard for me to understand any topic in this course. Since this alternative assessment is important as part of the courses, therefore I knew I will need to do some research any information regarding this project. Luckily, before starting the alternative assessment, me and my group members for project has started doing some project works especially in the text preprocessing step. Therefore, it wasn't awkward for me to start doing this alternative assessment by myself. I could say that I really enjoyed the experience of doing data mining task. I enjoyed doing the research on how to do data mining. The thing that I found frustrating about is that every knowledge that I applied during completing the alternative assessment, I could say that almost everything is wrong and misleading. That's when I realized that my efforts were kind of been wasted just like that. But nevertheless, it brings me to a new experience whereby I learnt from my mistakes.

(b) For me, the most difficult part in this alternative assessment is the part 1 question 2 where we needed to do data mining task. I found it's hard because I may can do the data mining task, but I don't really understand the result of it. That's when I knew I needed to do a lot more research than before. I gathered a lot more information about each supervised and unsupervised learning algorithm. The thing that I remembered to be effective is I started this alternative assessment earlier. As soon as Dr. Rozilawati gave it to us, I opened it and read everything through so that I could measure my strength of knowledge based on the difficulty of the questions. By the time I received this alternative assessment also, I knew that I could easily done the part 1 question 1. But the question 2, I seemed a bit taken aback since by that time my knowledge about doing data mining task is a bit low. The most ineffective I ever did during the completion of the alternative assessment is I didn't allocate the time wisely to complete the alternative assessment. After done with question 1 for part 1, I left the question 2 unanswered for a while because I thought that I needed to do more research about it. So, I did other courses' assignments first. Few days later, then I just realized that I haven't start anything for question 2, so I was panic and tried to do everything and completed the question 2, two days before the submission.

(c) There are few things that I think I needed to improve in this course. Firstly, in order for me to improve myself in this course, I should explore more other data mining tools. In the short time given, I chose RapidMiner for my tool for this alternative assessment because it has many resources for me to at least learn from. But I think after this I might want to explore tools such as Weka, Power BI etc. Secondly, I should expose myself with handling different type of dataset with a big amount of data. Since in the world, we have a lots of data type. Therefore, by having the experience to handle different types of data, it can enhance my knowledge about data mining, and I can also get the benefits from it. In my opinion, this Data Mining course can be applied in my future job. Since I'm major in Data Engineer, my scope of job definitely will be in this field too, which needed to handle big data. This course really helped to expose me to the real-world job especially in data.

(d) In a nutshell, in order for me to meet the learning needs, I should do a lot more findings about data mining task and engage myself handling big data more. To be success specifically in this course, I should try and error all data mining algorithm. Learn from mistakes and change to be better. Throughout this course, I found that I had developed new skills which are preprocess data from raw and process them through data mining task. Before entering this course, I don't even know how to clean data and I thought that I can't do it. After finishing this course, I realized that cleaning data wasn't really that hard if I practiced a lot. After all the practice and learning from our lecturer, Dr. Rozilawati, my skills increased

# Reference

1. *(2) What is the difference between clustering and association rule mining? - Quora.* (2018). Quora.com. https://www.quora.com/What-is-the-difference-between-clustering-and-association-rule-mining

2. Minewiskan. (2018, May 8). *Training and Testing Data Sets*. Microsoft.com. https://docs.microsoft.com/en-us/analysis-services/data-mining/training-and-testing-data-sets?view=asallproducts-allversions

3. https://www.facebook.com/MachineLearningMastery. (2018, May 22). *A Gentle Introduction to k-fold Cross-Validation*. Machine Learning Mastery. https://machinelearningmastery.com/k-fold-cross-validation/

4. *2.1 What is the difference between labelled and unlabelled data? · Grokking Machine Learning MEAP V14*. (2021). Manning.com. https://livebook.manning.com/book/grokking-machine-learning/2-1-what-is-the-difference-between-labelled-and-unlabelled-data-/v-4/40

# APPENDIX