

Clustering Based Approach to Enhance Association Rule Mining

MSc Research Project
MSc Data Analytics

Samruddhi Shailesh Kanhere

Student ID: x18190634

School of Computing
National College of Ireland

Supervisor: Dr. Paul Stynes, Dr. Pramod Pathak

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Samruddhi Shailesh Kanhere
Student ID:	x18190634
Programme:	MSc Data Analytics
Year:	2020
Module:	MSc Research Project
Supervisor:	Dr. Paul Stynes, Dr. Pramod Pathak
Submission Due Date:	28/09/2020
Project Title:	Clustering Based Approach to Enhance Association Rule Mining
Word Count:	1048
Page Count:	10

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	Samruddhi Shailesh Kanhere
Date:	28th September 2020

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project , both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Clustering Based Approach to Enhance Association Rule Mining

Samruddhi Shailesh Kanhere
x18190634

1 Introduction

The configuration manual is written so that the research work can be reproduced. The topic of my research is “**Clustering based approach to enhance association rule mining**”. It includes all the details about the System Setup i.e. hardware and software requirements, programming languages and libraries that are used. Further, this manual contains all the steps to run the code. In the last section, the important code snippets are attached.

2 System Setup

This section has the steps which will help in setting up the environment. The hardware and software requirements are listed below.

2.1 Hardware Requirements

The research is performed on a personal laptop. Below are the configurations of the same.

- Operating System: Windows OS, 64-bit
- Storage: 1 TB HDD
- RAM: 8.00 GB
- Processor: Intel(R) Core(TM) i5-8250U CPU @ 1.60GHz 1.80 GHz

2.2 Software Requirements

This section lists the softwares that needs to be installed before running the code.

- **Microsoft Excel 2016**

Excel is a widely used spreadsheet tool developed by Microsoft. This is used to store the data in the Comma Separated Values (CSV) format. This tool is also used for capturing the results of the experiments.

- **RStudio Desktop Version 1.2.1335**

RStudio is an open source software. It provides an Integrated Development Environment for R. It can be downloaded from RStudio official website¹. Before installing RStudio, R needs to be installed. R is an open source programming language specifically designed for statistical computing. It can be downloaded from CRAN website². The R version used in this research is 3.6.0. It is used for exploratory data analysis, data preparation, data modelling, etc.

- **Anaconda Navigator Distribution: Spyder 3.3.6**

Anaconda Navigator is an open source platform for Python operations. Python version 3.7 is used in this research. Python is used for Data modelling in this research. Anaconda can be downloaded from the official website of Anaconda³. Python can be downloaded from the Python's official website⁴.

- **Notepad ++ 7.7.1**

Notepad ++ is an open source editor. This is used for scratch work in this research. It can be downloaded from Notepad official website⁵.

3 Libraries

3.1 Python libraries

Python is mainly used for implementation of Frequent Pattern (FP) Growth algorithm. The Python libraries listed in Table 1 are used in this research. To install Python packages, use the following command in the anaconda prompt⁶.

conda install package_name

Table 1: Python libraries

Library	Version	Purpose
datetime	0.4-3	record the execution times
json	0.8.5	dealing with Json data
matplotlib	3.1.1	visualize the results
pandas	0.25.1	reading data
pyfpgrowth	1.0	implement FP Growth
seaborn	0.9.0	creating visualizations

¹<https://rstudio.com/products/rstudio/download/>

²<https://cran.r-project.org/bin/windows/>

³<https://www.anaconda.com/products/individual>

⁴<https://www.python.org/downloads/>

⁵<https://notepad-plus-plus.org/downloads/v7.7.1/>

⁶<https://docs.conda.io/projects/conda/en/latest/user-guide/tasks/manage-pkgs.html>

3.2 R libraries

The R libraries listed in Table 2 are used in this research. To install these R packages, following command should be executed⁷.

```
install.packages("package name")
```

Table 2: R libraries

Library	Version	Purpose
arules	1.6-6	implement association rule mining
arulesViz	1.3-3	visualize the association rules and frequent itemsets
clValid	0.6-9	validate the clustering results
clustree	0.4-3	visualizing the clusters
dplyr	1.0.0	data manipulation
factoextra	1.0.7	find the optimal number of clusters
FactoMineR	2.3	exploratory data analysis
ggfortify	0.4.10	Visualization of statistical results
GGally	2.0.0	creating plots
ggiraphExtra	0.2.9	making interactive plots
jsonlite	1.6	json parser
lubridate	1.7.4	date operations
NbClust	3.0	implementing clustering
plyr	1.8.4	tools to split and combine data
RColorBrewer	1.1-2	colour brewer palettes
tidyverse	1.2.1	text data manipulation

4 Code Snippets

This section includes the important code snippets from the overall project development process. The important steps in this process are data preparation, data modelling, and evaluation. Before that, Table 3 shows the important columns from the dataset and its description.

Table 3: Dataset Description

Column	Type	Description
RequestBasketValue	Decimal	Full basket value as per the bill
ResponseProcessingTimestamp	Date Time (UTC)	dealing with Json data
RequestBasketId	Varchar	Transaction id
RequestNumberBasketItems	Integer	Number of items in transaction
ResponseFinancialTimestamp	Date Time (UTC)	Transaction processing time
RequestBasketJsonString	Varchar	JSON string containing basket items

⁷http://jtlee.com/modules/01_DataScientistToolbox/02_09_installingRPackages/#11

4.1 Data Preparation

The first step in the project development process is data preparation. It includes all the operations related to data importing, data pre-processing and data transformation. These operations are performed in the RStudio.

```
#####Importing data#####
#Load the data files into the R environment
BigTempRb <- read.csv("E:/NCI/SEM3/RB/20181129_PMRB_0001.txt")
Products <- read.csv("E:/NCI/SEM2/RIC/RetailData/Promopay/tcg_products_20181219.csv")

#####Data Pre-processing#####
#Converting the JSON data into the structured format.
json <- BigTempRb$RequestBasketJsonString[1:300000] #Only taking 300,000 records due to system limitations
nested_dataframe <- jsonlite::stream_in(textConnection(gsub("\n", "", json)))
nested_dataframe <- cbind(nested_dataframe, dt=BigTempRb$ResponseFinancialTimestamp[1:300000], basket_items=BigTempRb$RequestNumberBasketItems[1:300000])
df <- nested_dataframe %>% unnest(items)
df <- df[-c(2)] # Remove the unwanted column Curr
names(Products)[1] <- "b" # Rename Product table EAN to b
df <- inner_join(df, Products, by="b")
df1 <- df

#Removing the items having negative quantity. Also, the baskets with single items are ignored.
df <- filter(df, b != '5000128785617')
df <- filter(df, b != '0000000000145')
df <- filter(df, basket_items != "1")
df <- filter(df, q >= 1)

head(df)
##### Changing Date Time format to Categorical #####
#For differential MBA, we need to change the time of transaction into a categorical variable.
#The transactions are classified into 4 classes: Morning, Afternoon, Evening, Night
library(lubridate)
library(plyr)
test1 <- mapvalues(hour(df$dt), from=c(0:23),
                  to=c(rep("Night",times=5), rep("Morning",times=6),rep("Afternoon",times=5),rep("Evening", times=4),rep("Night", times=4)))
df <- cbind(df, time_cat=test1)
head(df)
```

Figure 1: Data Preparation

As shown in Figure 1, the data pre-processing includes the conversion of semi-structured JSON data into structured format, handling missing values and inconsistencies, changing the date format, etc. This is how data is prepared for further modelling.

4.2 Data Modelling

Data modelling section is all about the implementation of the research. This section includes experiment wise code snippets.

- **Experiment 1: Replication of state-of-the-art (Hossain et al.; 2019).**

Two datasets are used for this experiment. These datasets are downloaded from Kaggle. This experiment is performed in Python language.

- **Dataset 1⁸: French Retail Dataset**

```
21
22 # reading the dataset
23 data = pd.read_csv('E:/NCI/SEM3/Thesis/Code/Dataset/StateOfArt/Market_Basket_Optimisation.csv', header = None)
24
25 #Data Pre-processing
26 transactions = []
27 for i in range(0, 7501):
28     transactions.append([str(data.values[i,j]) for j in range(0, 20) if str(data.values[i,j]) != 'nan'])
29
```

Figure 2: French Retail Dataset Preprocessing

Figure 2,3,4 shows the code snippets for french retail dataset.

- **Dataset 2⁹: Bakery dataset**

```

30
31 # Let's check the shape of the dataset
32 data.shape
33 plt.rcParams['figure.figsize'] = (15, 15)
34 wordcloud = WordCloud(background_color = 'white', width = 1200, height = 1200, max_words = 121).generate(str(data[0]))
35 plt.imshow(wordcloud)
36 plt.axis('off')
37 plt.title('Most Popular Items',fontsize = 20)
38 plt.show()
39
40 # Looking at the frequency of most popular items
41 plt.rcParams['figure.figsize'] = (18, 7)
42 color = plt.cm.copper(np.linspace(0, 1, 40))
43 data[0].value_counts().head(40).plot.bar(color = color)
44 plt.title('frequency of most popular items', fontsize = 20)
45 plt.xticks(rotation = 90 )
46 plt.grid()
47 plt.show()
48
49 y = data[0].value_counts().head(50).to_frame()
50 y.index
51
52 # plotting a tree map
53 plt.rcParams['figure.figsize'] = (20, 20)
54 color = plt.cm.cool(np.linspace(0, 1, 50))
55 squarify.plot(sizes = y.values, label = y.index, alpha=.8, color = color)
56 plt.title('Tree Map for Popular Items')
57 plt.axis('off')
58 plt.show()

```

Figure 3: French Retail Dataset: Exploratory Data Analysis

```

113
114 #Define threshold values for FP Growth algorithm
115 min_support = 0.01
116 min_confidence = 0.5
117 support = round(len(transactions)*min_support)
118
119 #Generating Frequent Itemsets
120 patterns = pyfpgrowth.find_frequent_patterns(transactions, support)
121
122 #Generating Association Rules
123 rules = pyfpgrowth.generate_association_rules(patterns, min_confidence)
124

```

Figure 4: French Retail Dataset: FP Growth Implementation

```

13 #Read the data
14 bakery_df = pd.read_csv('E:\NCI\SEM3\Thesis\Code\Dataset\StateOfArt\BakeryDataset.csv')
15 bakery_df.head()
16 bakery_df.info()
17
18 #Pre-processing the data
19 #Converting items to lower abd removing missing values.
20 bakery_df['Item'] = bakery_df['Item'].str.lower()
21 x = bakery_df['Item'] == "none"
22 print(x.value_counts())
23 bakery_df = bakery_df.drop(bakery_df[bakery_df.Item == 'none'].index)
24 len(bakery_df['Item'].unique())
25

```

Figure 5: Bakery Dataset Preprocessing

Figure 5,6,7 shows the code snippets for bakery dataset.

This experiment is performed several times and execution times are captured by changing the support values.

⁸<https://www.kaggle.com/roshansharma/market-basket-optimization>

⁹<https://www.kaggle.com/sulmansarwar/transactions-from-a-bakery>

```

25
26 #Top 20 top selling items
27 fig, ax=plt.subplots(figsize=(16,7))
28 bakery_df['Item'].value_counts().sort_values(ascending=False).head(20).plot.bar(width=0.5,edgecolor='k',align='center',linewidth=1)
29 plt.xlabel('Food Item',fontsize=20)
30 plt.ylabel('Number of transactions',fontsize=17)
31 ax.tick_params(labelsize=20)
32 plt.title('20 Most Sold Items at the Bakery',fontsize=20)
33 plt.grid()
34 plt.ioff()
35

```

Figure 6: Bakery Dataset: Exploratory Data Analysis

```

36 #Trasforming data into required format
37 grouped_df = bakery_df.groupby('Transaction')
38 grouped_lists = grouped_df["Item"].apply(list)
39 grouped_lists = grouped_lists.reset_index()
40
41 print(len(grouped_lists))
42 #Define threshold values
43 min_support = 0.01
44 min_confidence = 0.5
45 support = round(len(grouped_lists)*min_support)
46
47 #Calculating Frequent Itemsets
48 patterns = pyfpgrowth.find_frequent_patterns(grouped_lists["Item"], support)
49 #Generating Association Rules
50 rules = pyfpgrowth.generate_association_rules(patterns, min_confidence)
51 print(rules)
52

```

Figure 7: Bakery Dataset: FP Growth Implementation

- **Experiment 2: Implementing FP Growth for Glantus dataset.**

The FP Growth algorithm is implemented on Glantus dataset by keeping the same parameters as state-of-the-art (Hossain et al.; 2019). This experiment is performed in Python. Figure 8, 9 show the code snippet for Experiment 2.

```

14 ## Reading the transactions file
15 df=pd.read_csv("E:\\NCI\\SEM2\\RIC\\Code\\Dataset\\BigTempRb.txt",delimiter=',')
16
17 x = datetime.datetime.now()
18 print("Program beginning at :",x)
19 df1 = pd.DataFrame(columns=['ResponseRgBasketId', 'BasketID', 'ItemID','Qty','Price'])
20 ## json parser for parsing the json basket strings
21 itemlist=[]
22 for i in range(0,len(df[["ResponseRgBasketId"]])):
23     response=json.loads(df["RequestBasketJsonString"].iloc[i])
24     #for nest in response['items']:
25     lister = []
26     for nest in response['items']:
27         lister.append(nest["b"])
28     itemlist.append(list(dict.fromkeys(lister)))

```

Figure 8: Glantus Dataset: Data Preprocessing

- **Experiment 3: Performance Comparison of Apriori, FP Growth, and**


```

29
30 #Set same parameters as state-of-the-art
31 minimum_sup = 0.1
32 minimum_conf = 0.5
33 #FP Growth algorithm to generate the rules.
34 support = round(len(itemlist)*minimum_sup)
35 #Generate frequent patterns and association rules.
36 patterns = pyfpgrowth.find_frequent_patterns(itemlist, support)
37 rules = pyfpgrowth.generate_association_rules(patterns, 0.5)
38 print(rules)
~

```

Figure 9: Glantus Dataset: FP Growth Implementation

Eclat algorithms.

Glantus dataset is used for this experiment. These three association rule mining algorithms are implemented. The experiment is performed several times by to record the execution times by changing the number of input transactions in each iteration. The readings are recorded in Microsoft Excel.

– Apriori Algorithm

Apriori algorithm is implemented on the Glantus dataset after completing the pre-processing explained in Section 4.1. This is implemented in R language. Figure 10 shows the code snippet for the same.

```

75 #Track Time### Apriori Pre-processing#####
76 algo_start_time = Sys.time()
77 basket <- as(Transactions_by_Subcategory, "transactions")
78 itemFrequencyPlot(basket,topN=20,type="absolute",col=brewer.pal(8,'Pastel2'), main="Absolute Item Frequency Plot")
79 #itemFrequencyPlot(basket,topN=20,type="relative",col=brewer.pal(8,'Pastel2'),main="Relative Item Frequency Plot")
80 summary(basket)
81 ##### Apriori #####
82 association_rules <- apriori(basket, parameter = list(supp=0.002, conf=0.005,maxlen=10,target="rules",minlen=2))
83 summary(association_rules) #shows the following:
84 rules <- association_rules[!is.redundant(association_rules)]
85 as2221 <- inspectDT(sort(association_rules, by = "confidence"))
86 inspectDT(sort(association_rules, by = "confidence"))
87 algo_end_time = Sys.time()
88 algo_time = difftime(algo_end_time,algo_start_time,units = "secs")

```

Figure 10: Glantus Dataset: Apriori Algorithm

– Eclat Algorithm

In this section, Eclat algorithm is implemented using R language. Before doing that, data pre-processing is performed. Figure 11 shows the code snippet for the same.

```

90 ##### ECLAT Algorithm#####
91 eclat_frequent_itemsets = eclat(data = basket, parameter = list(support = 0.002, minlen = 2 ))
92 eclat_rules <- ruleInduction(eclat_frequent_itemsets, basket, confidence = 0.005)
93 inspectDT(sort(eclat_rules, by = 'support'))
94

```

Figure 11: Glantus Dataset: Eclat Algorithm

– FP Growth Algorithm

FP Growth algorithm is implemented using Python language. The values of input parameters are changed than the previous experiment. Figure 12 shows the code snippet for the same.

```

14 ## Reading the transactions file
15 my_df=pd.read_csv("E:\\NCI\\SEM3\\RB\\20181129_PMRB_0001.txt", delimiter=',')
16 df = my_df[1:100000]
17 x = datetime.datetime.now()
18 print("Program beginning at :",x)
19 df1 = pd.DataFrame(columns=['ResponseRgBasketId', 'BasketID', 'ItemID', 'Qty', 'Price'])
20
21 ## json parser for parsing the json basket strings
22 itemlist=[]
23 for i in range(0,len(df[["ResponseRgBasketId"]])):
24     response=json.loads(df["RequestBasketJsonString"].iloc[i])
25     #for nest in response['items']:
26     lister = []
27     for nest in response['items']:
28         lister.append(nest["b"])
29     itemlist.append(list(dict.fromkeys(lister)))
30
31 #Setting Parameters
32 minimum_sup = 0.002
33 minimum_conf = 0.005
34 #FP Growth algorithm to generate the rules.
35 support = round(len(itemlist)*minimum_sup)
36 patterns = pyfpgrowth.find_frequent_patterns(itemlist, support)
37 rules = pyfpgrowth.generate_association_rules(patterns, minimum_conf)
38 print(rules)
39 y = datetime.datetime.now()
40 print("Program ending at :",y)
41 print("total time : ", (y-x).total_seconds())

```

Figure 12: Glantus Dataset: FP Growth Algorithm

- **Experiment 4: Implement Clustering based approach for dataset reduction in Association Rule Mining.**

In this experiment, K-means clustering is implemented and based on the results of K-means the dataset is reduced. After dataset reduction, the association rule mining algorithms are implemented. The Figure 13 and Figure 14 show the code snippets for K-means clustering. The frequency and average price is calculated for all the products as a part of preprocessing. Then, K-means clustering is implemented based on these calculated parameters. The implementation is done in R. Figure 15 shows the snippet to filter the data based on clustering results.

- **Experiment 5: Implementing differential market basket analysis for Glantus dataset.**

Transactions are groups based on the time. The four groups are 'Afternoon', 'Morning', 'Night', and 'Evening'. Then, the association rule mining algorithms are implemented on each group and the output is compared. The implementation is performed in R. Figure 16 and Figure 17 show the code snippets for this experiment.

```

95- ##### K Means Clustering #####
96 Frequency_Table <- count(df, 'b')
97 library(dplyr)
98 #max_length_trans <- df %>% count(id)
99 #view(max_length_trans)
100 #view(Frequency_Table)
101 Frequency_Table$b[Frequency_Table$b == ""] <- NA
102 Frequency_Table <- na.omit(Frequency_Table, na.string = "")
103
104 test2 <- cbind(df$b, df$p)
105 test2 <- as.data.frame(test2)
106 test2[test2 == ""] <- NA
107 test2 <- na.omit(test2, na.string = "")
108
109 #Find the average price
110 test2$V2 <- as.numeric(as.character(test2$V2))
111 test2 <- aggregate(.~V1, data=test2, mean)
112 #Rename the columns
113 names(test2)[1] <- "b"
114 names(test2)[2] <- "price"
115 Frequency_Table <- inner_join(Frequency_Table, test2, by="b")
116

```

Figure 13: Glantus Dataset: Clustering Pre-processing

```

117 #Find the optimum numbers of clusters
118 fviz_nbclust(Frequency_Table, kmeans, method = "silhouette", k.max = 24) + theme_minimal() + ggtitle("The Silhouette Plot")
119
120 ## Optimum number of clusters is 4 after running "silhouette" Clustering for the data
121 set.seed(0)
122 km.res <- kmeans(Frequency_Table[,2:3], 4)
123
124 km.res$cluster
125 km.res$size
126 km.res$centers
127
128 # vary parameters for most readable graph
129 clusplot(Frequency_Table, km.res$cluster, color=TRUE, shade=TRUE,
130          labels=2, lines=0)

```

Figure 14: Glantus Dataset: Clustering Implementation

```

244- ##### Mapping cluster number to each transaction #####
245 cl_df <- inner_join(Frequency_Table[,c(1,4)], df, by="b")
246 cl_df$time_cat <- as.character(cl_df$time_cat)
247 cl_trans <- cl_df
248 #view(cl_trans)
249 cl_trans <- filter(cl_trans, cluster != '2')
250 a <- cl_trans

```

Figure 15: Filtering data based on Clustering

```

143 #####clustering the data based on the time of the day
144 df_morning <- cl_df %>% filter (cl_df$time_cat=='Morning')
145 df_afternoon <- cl_df %>% filter (cl_df$time_cat=='Afternoon')
146 df_evening <- cl_df %>% filter (cl_df$time_cat=='Evening')
147 df_night <- cl_df %>% filter (cl_df$time_cat=='Night')
148

```

Figure 16: Grouping data based on time

Figure 17 shows the Apriori algorithm implemented on Morning transactions. Similarly, it is implemented on the remaining groups.

```

150. ##Calculating rules for morning time.#####
151 df_morning$id <- as.factor(df_morning$id)
152 Morning_Transactions_by_Subcategory <- split(df_morning$b, df_morning$id)
153 #preproc_end_time = Sys.time()
154 #preproc_time = difftime(preproc_end_time,preproc_start_time, units = "secs")
155
156 #Track Time### Apriori Pre-processing
157 algo_start_time = Sys.time()
158 morning_basket <- as(Morning_Transactions_by_Subcategory, "transactions")
159 summary(morning_basket)
160 ##### Apriori #####
161 morning_association_rules <- apriori(morning_basket, parameter = list(supp=0.002, conf=0.005,maxlen=10,target="rules",minlen=2))
162 summary(morning_association_rules) #shows the following:
163 morning_association_rules <- morning_association_rules[!is.redundant(morning_association_rules)]
164 inspectDT(sort(morning_association_rules, by = "confidence"))
165 algo_end_time = Sys.time()
166 algo_time = difftime(algo_end_time,algo_start_time,units = "secs")
167

```

Figure 17: Apriori on Morning group

References

Hossain, M., Sattar, S. A. H. M. and Paul, M. K. (2019). Market Basket Analysis Using Apriori and FP Growth Algorithm, *2019 22nd International Conference on Computer and Information Technology (ICCIT)* pp. 1–6.