# D212 – OFM3 TASK 3 – ASSOCIATION RULES AND LIFT ANALYSIS

**Part I: Research Question**

A. Describe the purpose of this data mining report by doing the following:
1. Propose **one** question relevant to a real-world organizational situation that you will answer using market basket analysis.

What are the products that are usually bought together? When stakeholders have this knowledge, they will be able to come up with a marketing strategy to offer these products for example, at a discounted price, possibly reducing the customer churn.

2. Define **one** goal of the data analysis. Ensure that your goal is reasonable within the scope of the scenario and is represented in the available data.

Since acquiring a new customer is 10 times more costly than keeping one, the number one goal is to retain customers. In order to achieve this, stakeholders have to understand not only customer behavior in order to predict churn, but also what kind of products are interesting to a customer. To lower churn probability, stakeholders would be able to offer products that interest customers found in the same rule.

**Part II: Market Basket Justification**

B. Explain the reasons for using market basket analysis by doing the following:
1. Explain how market basket analyzes the selected dataset. Include expected outcomes.

Market Basket Analysis is a data mining technique that identifies products that exhibit strong relationships[1]. Using this analysis we will find out which combination of products are usually bought together (by association rule) and what is the relationship they present.

An expected outcome is to find the optimal combination of items to offer in a possible discounted price.

2. Provide **one** example of transactions in the dataset.

Every row represents a customer. Columns indicate what each customer bought (a customer can buy from 1 item to 20 items). For example, our second customer (row 5), bought 3 items:
- Apple Lightning to Digital AV Adapter;
- TP-Link AC1750 Smart WiFi Router and
- Apple Pencil.

3. Summarize **one** assumption of market basket analysis.

One assumption in market basket analysis is when two or more products are bought together, they are complementary, and therefore purchase of one will lead to purchase of the other(s).

**Part III: Data Preparation and Analysis**
C. Prepare and perform market basket analysis by doing the following:
1. Transform the dataset to make it suitable for market basket analysis. Include a copy of the cleaned dataset.

In order to make the provided dataset suitable for market basket analysis, first we need to clean the dataset then we have to create lists of the itemsets. The cleaned dataset will be uploaded along with this document.

Cleaning the dataset:

```python
#Finding missing values in my dataset
churn_df.isnull().any(axis=1)
null_values = churn_df.isna().any()
print(null_values)

#Lets drop missing values
churn_df.dropna(how='all', inplace = True)

churn_df.fillna(0, inplace = True)

#New Dataset
print('Clean Dataset Shape: ', churn_df.shape)

#Checking if we have "0" in the dataset
print(churn_df.head())

churn_df.info()
```

Converting to list:

```python
#Converting to list so we can use Apriori Algorithm
churn_df_list = []
for i in range(0, 7501):
    churn_df_list.append([str(churn_df.values[i,j])for j in range(0,20)])
churn_df_clean = pd.DataFrame(churn_df_list)

print(churn_df_clean.head())
```

Extracting the clean dataset:

```python
#Extract the "Prepared"" dataset
churn_df_clean.to_csv('prepared_churn_data_mba.csv')
churn_df = pd.read_csv('prepared_churn_data_mba.csv')
df = churn_df.columns
print('The dataset columns are ', df)
```

**Figure 1: Screen Shot of Dataset for MBA**

2. Execute the code used to generate association rules with the Apriori algorithm. Provide screenshots that demonstrate the error-free functionality of the code.

The **apriori algorithm** is an efficient alternative that helps identify frequent itemsets while filtering out the infrequent ones. It can remove itemsets from consideration without having to evaluate them[1].

```
#Training the Algorithm
apriori_list = apriori(churn_df_list, min_support=0.003, min_confidence=0.3,
min_lift=3, min_lenght=2)

apriori_list = list(apriori_list)
print(apriori_list[0])

dataset = pd.DataFrame(apriori_list)

print(dataset)

#Since I couldnt visualize all columns in dataset (items, support,
ordered_stats)
print(dataset.columns)

#Lets separate columns in their own
support = dataset.support

#Antecedent, Consequence, confidence, lift
antecedent_values = []
consequent_values = []
confidence_values = []
lift_values = []

for i in range(dataset.shape[0]):
    single_list = dataset['ordered_statistics'][i][0]
    antecedent_values.append(list(single_list[0]))
    consequent_values.append(list(single_list[1]))
```

```python
        confidence_values.append(single_list[2])
        lift_values.append(single_list[3])

#From List to Dataframe
antecedent = pd.DataFrame(antecedent_values)
consequent = pd.DataFrame(consequent_values)
#The confidence metric measures the likelihood of a consequent being
purchased given the purchase of the antecedent.
confidence = pd.DataFrame(confidence_values, columns=['CONFIDENCE'])
#The lift metric measures the influence that the purchase of the antecedent
has on the purchase of the consequent.
lift = pd.DataFrame(lift_values, columns=['LIFT'])

#Concatenate lists into a dataframe
dataframe = pd.concat([antecedent, consequent, support, confidence, lift],
axis = 1)
dataframe.fillna(value = '', inplace= True)

print(dataframe)

print(dataframe.columns)
print(dataframe.head())
```
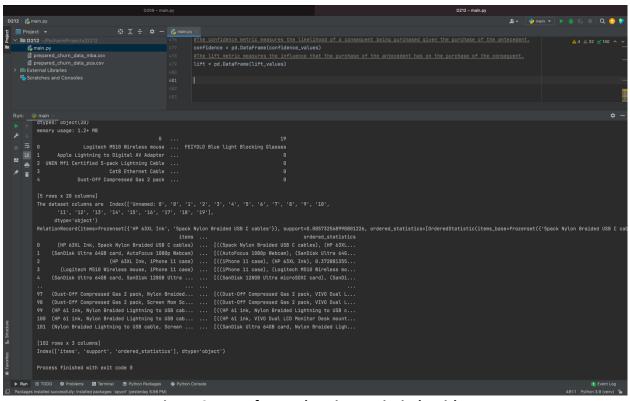


**Figure 2: Error free code using Apriori Algorithm**

3. Provide values for the support, lift, and confidence of the association rules table.

```
#Association Rules
dataframe.columns = ['antecedent', 1, 2, 'consequent', 1, 2, 'support',
'CONFIDENCE', 'LIFT']
print(dataframe.columns)
dataframe_1 = dataframe[['antecedent', 'consequent', 'support', 'CONFIDENCE',
'LIFT']]

#Export the dataframe_1 into csv
dataframe_1.to_csv('dataset_mba_association.csv')
```

Attached to this document it will be also uploaded another spreadsheet called dataset_mba_association.csv. This spreadsheet contains all support, lift and confidence values for the association rules table. Here is a screen shot of it:

| | antecedent | consequent | support | CONFIDENCE | LIFT | |
|---|---|---|---|---|---|---|
| 0 | 5pack Nylon Braided USB C cables | HP 63XL Ink | 0.005732569 | 0.300699301 | 3.790832697 | |
| 1 | AutoFocus 1080p Webcam | SanDisk Ultra 64GB card | 0.005332622 | 0.377358491 | 3.840659481 | |
| 2 | iPhone 11 case | HP 63XL Ink | 0.005865885 | 0.372881356 | 4.70081185 | |
| 3 | iPhone 11 case | Logitech M510 Wireless mouse | 0.005065991 | 0.322033898 | 4.506672148 | |
| 4 | SanDisk 128GB Ultra microSDXC card | SanDisk Ultra 64GB card | 0.015997867 | 0.323450135 | 3.291993841 | |
| 5 | 5pack Nylon Braided USB C cables | HP 63XL Ink | 0.005732569 | 0.300699301 | 3.790832697 | |
| 6 | AutoFocus 1080p Webcam | SanDisk Ultra 64GB card | 0.005332622 | 0.377358491 | 3.840659481 | |
| 7 | iPhone 11 case | HP 63XL Ink | 0.005865885 | 0.372881356 | 4.70081185 | |
| 8 | iPhone 11 case | Logitech M510 Wireless mouse | 0.005065991 | 0.322033898 | 4.515095834 | |
| 9 | SanDisk 128GB Ultra microSDXC card | SanDisk Ultra 64GB card | 0.015997867 | 0.323450135 | 3.291993841 | |
| 10 | Anker USB C to HDMI Adapter | Screen Mom Screen Cleaner kit | 0.003066258 | 0.442307692 | 3.413323045 | |
| 11 | FEIYOLD Blue light Blocking Glasses | Screen Mom Screen Cleaner kit | 0.00359952 | 0.402985075 | 3.10986733 | |
| 12 | HP 61 ink | Screen Mom Screen Cleaner kit | 0.003999467 | 0.394736842 | 3.046215075 | |
| 13 | Nylon Braided Lightning to USB cable | Screen Mom Screen Cleaner kit | 0.003999467 | 0.5 | 3.858539095 | |
| 14 | SanDisk Ultra 64GB card | Screen Mom Screen Cleaner kit | 0.003999467 | 0.410958904 | 3.171401996 | |
| 15 | Dust-Off Compressed Gas 2 pack | FEIYOLD Blue light Blocking Glasses | 0.003866151 | 0.402777778 | 6.115862573 | |
| 16 | VIVO Dual LCD Monitor Desk mount | Screen Mom Screen Cleaner kit | 0.003999467 | 0.454545455 | 3.507762813 | |
| 17 | Anker USB C to HDMI Adapter | VIVO Dual LCD Monitor Desk mount | 0.004399413 | 0.611111111 | 3.509911519 | |
| 18 | Anker USB C to HDMI Adapter | Nylon Braided Lightning to USB cable | 0.003999467 | 0.357142857 | 3.746753247 | |
| 19 | Anker USB C to HDMI Adapter | Nylon Braided Lightning to USB cable | 0.003066258 | 0.365079365 | 3.83001443 | |
| 20 | Anker USB C to HDMI Adapter | Nylon Braided Lightning to USB cable | 0.006665778 | 0.318471338 | 3.341053851 | |
| 21 | Apple Pencil | SanDisk Ultra 64GB card | 0.004132782 | 0.329787234 | 3.356491238 | |
| 22 | Apple Pencil | VIVO Dual LCD Monitor Desk mount | 0.003732836 | 0.528301887 | 3.034297437 | |
| 23 | AutoFocus 1080p Webcam | VIVO Dual LCD Monitor Desk mount | 0.003066258 | 0.575 | 3.302507657 | |
| 24 | Logitech M510 Wireless mouse | Nylon Braided Lightning to USB cable | 0.00719904 | 0.305084746 | 3.200616333 | |
| 25 | SanDisk 128GB Ultra microSDXC card | SanDisk Ultra 64GB card | 0.006665778 | 0.390625 | 3.975682666 | |
| 26 | FEIYOLD Blue light Blocking Glasses | Screen Mom Screen Cleaner kit | 0.00359952 | 0.5 | 3.858539095 | |
| 27 | Logitech M510 Wireless mouse | Screen Mom Screen Cleaner kit | 0.003199573 | 0.393442623 | 3.036227484 | |
| 28 | FEIYOLD Blue light Blocking Glasses | Screen Mom Screen Cleaner kit | 0.00479936 | 0.423529412 | 3.268409586 | |
| 29 | Falcon Dust Off Compressed Gas | Screen Mom Screen Cleaner kit | 0.003866151 | 0.408450704 | 3.152046021 | |
| 30 | HP 61 ink | SanDisk Ultra 64GB card | 0.003999467 | 0.441176471 | 4.490182776 | |
| 31 | HP 65 ink | Screen Mom Screen Cleaner kit | 0.003332889 | 0.416666667 | 3.215449246 | |
| 32 | Logitech M510 Wireless mouse | VIVO Dual LCD Monitor Desk mount | 0.0059992 | 0.523255814 | 3.00531536 | |
| 33 | VIVO Dual LCD Monitor Desk mount | SanDisk Ultra 64GB card | 0.008665511 | 0.311004785 | 3.165328209 | |
| 34 | SanDisk Ultra 64GB card | VIVO Dual LCD Monitor Desk mount | 0.00479936 | 0.571428571 | 3.281995187 | |
| 35 | VIVO Dual LCD Monitor Desk mount | SanDisk Ultra 64GB card | 0.005332622 | 0.322580645 | 3.283144395 | |
| 36 | Screen Mom Screen Cleaner kit | SanDisk Ultra 64GB card | 0.00359952 | 0.391304348 | 3.982596897 | |
| 37 | USB 2.0 Printer cable | SanDisk Ultra 64GB card | 0.003199573 | 0.461538462 | 4.697421981 | |
| 38 | VIVO Dual LCD Monitor Desk mount | SanDisk Ultra 64GB card | 0.006399147 | 0.393442623 | 4.004359722 | |
| 39 | SanDisk Ultra 64GB card | VIVO Dual LCD Monitor Desk mount | 0.003066258 | 0.676470588 | 3.885303126 | |
| 40 | VIVO Dual LCD Monitor Desk mount | SanDisk Ultra 64GB card | 0.003332889 | 0.337837838 | 3.438428252 | |
| 41 | Anker USB C to HDMI Adapter | Screen Mom Screen Cleaner kit | 0.003066258 | 0.442307692 | 3.413323045 | |
| 42 | FEIYOLD Blue light Blocking Glasses | Screen Mom Screen Cleaner kit | 0.00359952 | 0.402985075 | 3.10986733 | |
| 43 | HP 61 ink | Screen Mom Screen Cleaner kit | 0.003999467 | 0.394736842 | 3.046215075 | |
| 44 | Nylon Braided Lightning to USB cable | Screen Mom Screen Cleaner kit | 0.003999467 | 0.5 | 3.858539095 | |
| 45 | SanDisk Ultra 64GB card | Screen Mom Screen Cleaner kit | 0.003999467 | 0.410958904 | 3.171401996 | |
| 46 | Dust-Off Compressed Gas 2 pack | FEIYOLD Blue light Blocking Glasses | 0.003866151 | 0.402777778 | 6.128267974 | |
| 47 | VIVO Dual LCD Monitor Desk mount | Screen Mom Screen Cleaner kit | 0.003999467 | 0.454545455 | 3.507762813 | |

dataset_mba_association      +

**Figure 3: Association Rules Results**

4. Identify the top **three** rules generated by the Apriori algorithm. Include a screenshot of the top rules along with their summaries.

Based on the set parameters initially defined:

```
#Training the Algorithm
apriori_list = apriori(churn_df_list, min_support=0.003, min_confidence=0.3,
min_lift=3, min_lenght=2)
```

For these set of parameters, the algorithm returned 101 rules. Let's analyze the top three rules generated by the algorithm.

Rule 1 is the most relevant rule that the algorithm identified from the given dataset:
**Antecedent: 5pack Nylon Braided USB C cables**
**Consequent: HP 63XL Ink**
**Support: 0.00573256899080122**
**Confidence: 0.3006993006993**
**Lift: 3.79083269671504**

From all customers who purchased the cables, only 30% also purchased the printer ink. Support means that from all transactions, only 0.57% contain both items. The lift value means that once customers have purchased the cables, they have 3.8 times more chances to also purchase the printer's ink.

Rule 2:
**Antecedent: AutoFocus 1080p Webcam**
**Consequent: SanDisk Ultra 64GB card**
**Support: 0.00533262231702439**
**Confidence: 0.377358490566037**
**Lift: 3.84065948132408**

In this case here, 37% of all customers who purchased the webcam also purchased the memory card. From all transactions, only 0.533% contain both items and customers who buy the webcam have 3.8 more chances in also buying the memory card.

Rule 3:
**Antecedent: iPhone 11 case**
**Consequent: HP 63XL Ink**
**Support: 0.00586588454872683**
**Confidence: 0.372881355932203**
**Lift: 4.70081185016379**

In the last case, 37% of all customers who purchased the iPhone case also bought the ink. From all transactions, only 0.586% contain both of these items and when customers buy the phone case, they have 4.7 more likely in also acquiring the ink.

**Part IV: Data Summary and Implications**
D. Summarize your data analysis by doing the following:
   1. Summarize the significance of support, lift, and confidence from the results of the analysis.

Each pair of products A and B is evaluated on three measures[2]:

Support—the joint probability of finding the pair AB across all baskets. A low support means that the pair is not relevant because it is not purchased frequently enough.

Confidence—the conditional probability $p(B|A) = p(A \cap B)/p(A)$, which is often interpreted as the probability that purchase of product A will lead to purchase of product B.

Lift - measures the influence that the purchase of the antecedent has on the purchase of the consequent: $Lift(A,B) = \frac{Support(A,B)}{Support(A)*Support(B)}$. From this formula, we are looking for values greater than 1, meaning that the purchase of the antecedent increases the likelihood of the purchase of the consequent[1].

The results of my analysis did not show a good level of confidence. None of the 3 most relevant rules has a confidence level higher than 40%, far away from any good level.

We also obtained a very low support number, 0.5% which indicates that for the most relevant 3 rules, the purchase of the pair AB doesn't occur in more than half a percent.

Our highest lift is "4.7 more likely in buying" in the third rule, which means that when a customer buys the antecedent product "iPhone11 case" the same customer is 4.7 more likely to acquire the precedent product "XP 63XL ink".

From these results we can say that our analysis did not achieve good measures of support, confidence and lift.

2. Discuss the practical significance of the findings from the analysis.

All measures are very low and I would say we do not have a practical significance of the findings from the analysis. Perhaps more data needs to be gathered and we need to reevaluate our data one more time.

3. Recommend a course of action for the real-world organizational situation from part A1 based on your results from part D1.

Our analysis did not bring any significant result. None of the 3 most relevant pairings would suggest any difference in the current marketing/customer service strategy for the telecom company.

**Part V: Attachments**
E. Provide a Panopto video recording that includes a demonstration of the functionality of the code used for the analysis and a summary of the programming environment.

Video Link: https://wgu.hosted.panopto.com/Panopto/Pages/Viewer.aspx?id=fc245e35-ed14-4440-8a9b-ae0e01446d10

F.  Record *all* web sources used to acquire data or segments of third-party code to support the application. Ensure the web sources are reliable.

G.  Acknowledge sources, using in-text citations and references, for content that is quoted, paraphrased, or summarized.

[1] (2021, Nov 15th) NAIR, Aashish , Understanding Consumer Behavior With The Basket Market Analysis https://towardsdatascience.com/understanding-consumer-behavior-with-the-market-basket-analysis-3d0c017e5613

[2] (2012, May 22nd) KAMAKURA, Wagner, Sequencial Market Basket Analysis http://wak2.web.rice.edu/bio/My%20Reprints/Sequential%20Market%20Basket%20Analysist.pdf