# PROJECT REPORT

Supermarket Selles Transction Data Mining and Data Analysis Using Association Rule

Reporter: Eyasu Taye, Kalkidan Tesfaye, Chala Bahiru, Bilisie Melese

Department: Computer Science

# CONTENTS:

Part I: Introduction

Part II: Motivation

Part III: Objective

Part IV: Methodology

Part V: Methodology



## Introduction





Supermarkets have large number of customers checking into the items and to know customers need its better to identify which products bought frequently, which products are bought together, and association between items.

Association rule mining is one of the principal problems treated in KDD and can be defined as extracting the interesting correlation and relation among huge number of transactions.

Frequent itemset is generally adopted to generate association rules. As the amount of data stored supermarket database grows twice as fast as the speed of the fastest processor available to analyse it. Main purpose of analysing frequent itemset is to find the association relationship among the large number of database items. It is used to describe the patterns of customers' purchase in the supermarket.



### **Motivation**

Our motivation behined mining the supermarket data?

Finding inherent regularities in data is the motivation behind this supermarket data analysis project. Association rule mining is one of the technique to identify underlying relations between different items, it helps to identify which items of supermarket mostly bought together and their correlations.

# For instance, if item A and B are bought together more frequently then several steps can be taken to increase the profit.

### 1.Do no go far away

A and B can be placed together so that when a customer buys one of the product he/she doesn't have to go far away to buy the other product.

#### 3. Collective Discount

Collective discounts can be offered on those products, if the customer buys both.

### 2.Adverisement Campaign.

People who buy one of the products can be targeted through an advertisement campaign to buy the other one.

### 4. What to buy together

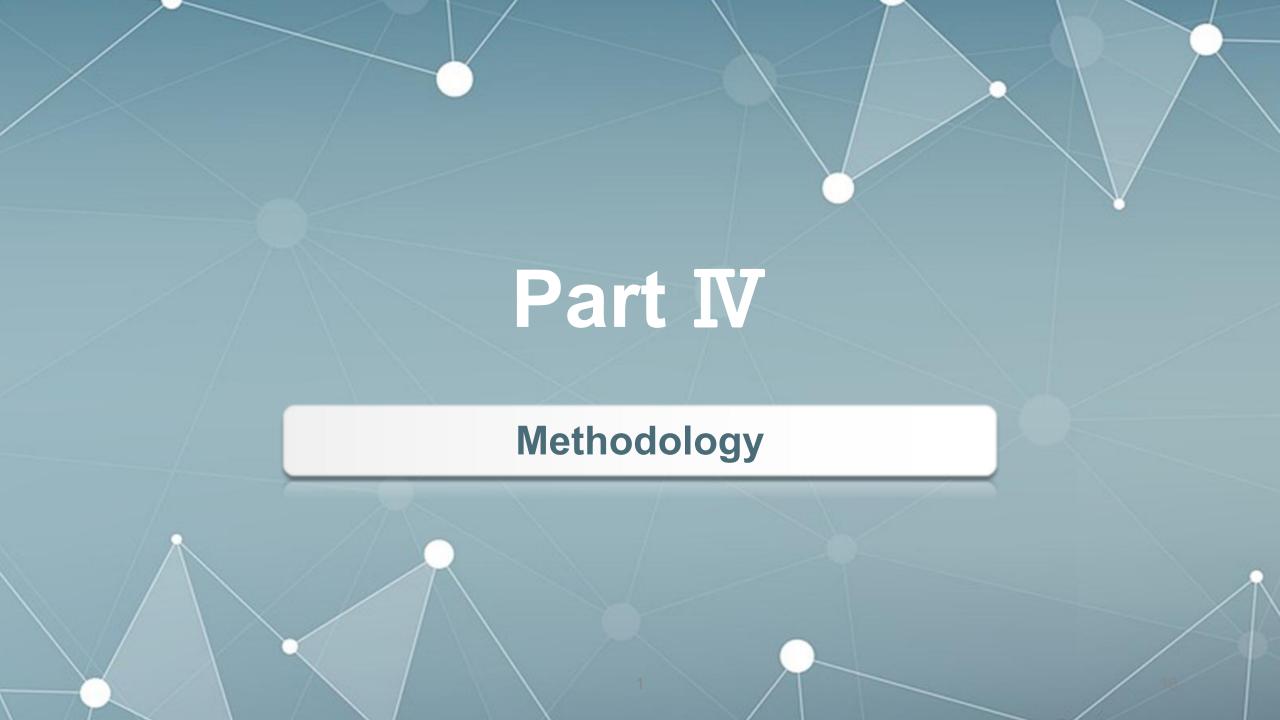
Both A and B can be packaged together.



# Objective

### What do we want to achieve?

The objective of analysing supermarket transaction data is for identification of items that frequently occurred together in the traniction found in the databases so that the out come will increase effectiveness in supermarket sell.



# Methodology

## **Data Preprocessing**

Step 1

Step 2

Step 3

The dataset has the following limitations

- Incomplete and Null values
- No header row is necessary
- Unused column

### **Data Cleaning**

- Make complete the data by inserting default value for null values.
- 2. Remove header option from the dataset since it is not necessary so the first row of the dataset can not be treated as header.

### **Data Reduction**

Remove unused column from the dataset.

### **Data Trnsformation**

Transform the csv dataset to list data type.

# Methodology



So the data is become understandable by the machine to be processed.

That is

Each row corresponds to a transaction and,

Each column corresponds to an item purchased in that specific transaction.

### The dataset dimension

```
In [49]: #check its dimension store_data.shape
```

Out[49]: (7500, 27)

7500 rows 27 columns

## Columns which are no more important in the dataset

t[89]:	shrimp	Ø		
1	almonds	1754		
	avocado	3112	2.2	9_202925
	vegetables mix	4156	olive oil	7500
	green grapes	4972	meat	7500
	whole weat flour yams cottage cheese energy drink	5637		
		6132 652 <b>0</b>	onion Are not importa	ant for 7500
		6847		
	tomato juice	7106	garlic our objectives	7500
	low fat yogurt	7245	dairy	7500
	green tea	7347		
	honey	7414	apples Must be remov	red 7500
	salad	rer 7476 seafood 7493 rejuice 7497 bananas othie 7497		
	mineral water salmon		searood	7500
	antioxydant juice		hananas	7500
	frozen smoothie			7500
	spinach	7498	dtype: int64	
_	olive oil	7500		
	meat	7500		
	onion	7500		
	garlic	7500 7500		
	dairy apples	7500		
	seafood	7500 7500		
	bananas	7500		
	dtype: int64		1	4.

### The first row as a header

# The first row is treated as a header, must be transformed to normal row

### True = The value is null False = The value is not null

•••	transfermed to morniar for											-	aioo	-		Talas it			
	shrimp	almonds	avocado	vegetables mix	green grapes	whole weat flour	yams	cottage cheese	energy drink	to mato juice		frozen smoothie	spinach	olive oil	meat	pasta	eo ups	dairy	egį
0	False	False	False	True	True	True	True	True	True	True		True	True	True	True	True	True	True	Tri
1	False	True	True	True	True	True	True	True	True	True	***	True	True	True	True	True	True	True	Tri
2	False	False	True	True	True	True	True	True	True	True		True	True	True	True	True	True	True	Tri
3	False	False	False	False	False	True	True	True	True	True	100	True	True	True	True	True	True	True	Tri
4	False	True	True	True	True	True	True	True	True	True		True	True	True	True	True	True	True	Tri
-	***		1000	225	22		33	(33)		(60)	321		125			125	3		
7495	False	False	False	True	True	True	True	True	True	True		True	True	True	True	True	True	True	Tri
7496	False	False	False	False	False	False	True	True	True	True	œ	True	True	True	True	True	True	True	Trı
7497	False	True	True	True	True	True	True	True	True	True	:::	True	True	True	True	True	True	True	Tri
7498	False	False	True	True	True	True	True	True	True	True	100	True	True	True	True	True	True	True	Tri
7499	False	False	False	False	True	True	True	True	True	True	12.2	True	True	True	True	True	True	True	Tri

### Data reduction by removing irrelevant columns

```
In [57]: #as we can see from the result the 20th to 27th column has all null values so we can elimintate it.
        store data.drop(store data.columns[[19,20,21,22,23,24,25,26]], axis=1, inplace=True)
In [58]: store_data.isnull().sum()
Out[58]: 0
               0
            1754
                                                 olive oil
            3112
                                                                                    7500
            4156
                                                 meat
                                                                                    7500
            4972
                                                 onion
            5637
                                                                                    7500
            6132
                                                 garlic
                                                               8 Columns
                                                                                    7500
            6520
             6847
                                                 dairy
                                                                                    7500
                                                                    are
            7106
                                                 apples
                                                                                    7500
       10
            7245
                                                                Removed
       11
            7347
                                                 seafood
                                                                                    7500
       12
            7414
                                                 bananas
                                                                                    7500
       13
            7454
       14
                                                 dtype: int64
            7476
       15
            7493
       16
            7497
       17
             7497
        18
            7498
        dtype: int64
```

### A header row is removed Null is transformed to NaN value

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
shrimp	almonds	avocado	vegetables mix	green grapes	whole weat flour	yams	cottage cheese	energy drink	tomato juice	low fat yogurt	green tea	honey	salad	mineral water	salmon	antioxydant juice	frozen smoothie	spinach
ourgers	meatballs	eggs	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
nutney	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
turkey	avocado	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
mineral water	milk	energy bar	whole wheat rice	green tea	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

A header row is removed Null is transformed to NaN value

<sup>\*</sup> Now the dataset is clean enough

# **Data Analysis**

### The apriori class requires some parameter values to work.

- The first parameter must be list of list that to extract rules from.
- The second parameter is the min support parameter.
- The min confidence parameter.



We want to create rules for items purchased five times a day in one week

That means

7\*5 = 35 times per week,

So the support = 35/7500 = 0.0045



Let

by this parameters 48 rules are dicovered

- he min\_lift parameter specifies the minimum lift value for the short listed rules.
- Finally, the min length parameter specifies the minimum number of items that you want in your rules.
- Minimum confidence = 20% or 0.2
- Lift = 3
- Min length = 2



### **Evalution**

#### **First Rule**

Rule: chicken -> light cream Support: 0.004532728969470737 Confidence: 0.29059829059829057 Lift: 4.84395061728395

- The support value is 0.0045. This number is calculated by dividing the number of transactions containing light cream divided by total number of transactions.
- The support value is 0.0045. This number is calculated by dividing the number of transactions containing light cream divided by total number of transactions.
- The confidence level for the rule is 0.2905 which shows that out of all the transactions that contain light cream, 29.05% of the transactions also contain chicken.
- The lift of 4.84 tells us that chicken is 4.84 times more likely to be bought by the customers who buy light cream compared to the default likelihood of the sale of chicken.

### **Evalution**

#### **Second Rule**

Rule: mushroom cream sauce -> escalope Support: 0.005732568990801226 Confidence: 0.3006993006993007 Lift: 3.790832696715049

- The second rule states that mushroom cream sauce and escalope are bought frequently
- The support for mushroom cream sauce is 0.0057.
- The confidence for this rule is 0.3006 which means that out of all the transactions containing mushroom, 30.06% of the transactions are likely to contain escalope as well.
- Lift of 3.79 shows that the escalope is 3.79 more likely to be bought by the customers that buy mushroom cream sauce, compared to its default sale.

## Conclusion



Association rule mining for supermarket dataset has been presented, Mining has been applied to sales data of dataset. In proposed project, the apriori algorithm has been used on super market dataset which gives associations of two products which has maximum support, It reduces the size of the itemsets in the database considerably providing a good performance. Thus, data mining helps consumers and industries better in the decision-making process.

# Thank you for watching

Click here to add the text, the text is the extraction of your thought, in order to finally present the good effect of the release, please try to be concise and concise; if necessary, add or subtract the text.

Reporter:XXXX

Department:XXXXX

```
In [1]: # Importing libraries
In [2]: import numpy as np
    import matplotlib.pyplot as plt
    import pandas as pd
    from apyori import apriori

In [3]: # Loading the dataset

In [4]: store_data = pd.read_csv('I:\\Datasets\\store_data.csv')

In [5]: #checking the dataset

In [6]: #check its dimension
    store_data.shape

Out[6]: (7500, 27)

In [7]: #check null values
    store_data.isnull()
```

Out[7]:

	shrimp	almonds	avocado	vegetables mix	green grapes	wnoie weat flour	yams	cottage cheese	energy drink	tomato juice	
0	False	False	False	True	True	True	True	True	True	True	
1	False	True	True	True	True	True	True	True	True	True	•••
2	False	False	True	True	True	True	True	True	True	True	
3	False	False	False	False	False	True	True	True	True	True	
4	False	True	True	True	True	True	True	True	True	True	•••
											•••
7495	False	False	False	True	True	True	True	True	True	True	•••
7496	False	False	False	False	False	False	True	True	True	True	•••
7497	False	True	True	True	True	True	True	True	True	True	•••
7498	False	False	True	True	True	True	True	True	True	True	•••
7499	False	False	False	False	True	True	True	True	True	True	

7500 rows × 27 columns

```
In [8]:
           #count null values and convert tehm to NaN default value
          store_data.isnull().sum()
 Out[8]: shrimp
          almonds
                                  1754
          avocado
                                  3112
          vegetables mix
                                  4156
          green grapes
                                  4972
          whole weat flour
                                  5637
          vams
                                  6132
          cottage cheese
                                  6520
          energy drink
                                  6847
          tomato juice
                                  7106
          low fat yogurt
                                  7245
          green tea
                                  7347
          honey
                                  7414
                                  7454
          salad
          mineral water
                                  7476
          salmon
                                  7493
          antioxydant juice
                                  7497
          frozen smoothie
                                  7497
          spinach
                                  7498
          olive oil
                                  7500
          meat
                                  7500
          onion
                                  7500
          garlic
                                  7500
          dairy
                                  7500
          apples
                                  7500
          seafood
                                  7500
          bananas
                                  7500
          dtype: int64
 In [9]:
          #for our processing we do not need a header row
          store_data = pd.read_csv('I:\\Datasets\\store_data.csv', header=None)
In [10]:
In [11]:
          store_data.head()
Out[11]:
                   0
                            1
                                    2
                                              3
                                                     4
                                                           5
                                                                 6
                                                                        7
                                                                                8
                                                                                      9 ...
                                                        whole
                                       vegetables
                                                  green
                                                                    cottage
                                                                           energy tomato
                      almonds avocado
           0
              shrimp
                                                         weat yams
                                            mix grapes
                                                                    cheese
                                                                             drink
                                                                                    juice
                                                                                            smo
                                                         flour
              burgers meatballs
                                            NaN
                                                   NaN
                                                         NaN
                                                               NaN
                                                                      NaN
                                                                             NaN
                                 eggs
                                                                                    NaN
             chutney
                         NaN
                                 NaN
                                            NaN
                                                   NaN
                                                         NaN
                                                               NaN
                                                                      NaN
                                                                             NaN
           2
                                                                                    NaN
               turkey
                      avocado
                                 NaN
                                            NaN
                                                   NaN
                                                         NaN
                                                               NaN
                                                                      NaN
                                                                             NaN
                                                                                    NaN
           3
                                          whole
              mineral
                                energy
                                                  green
                          milk
                                                         NaN
                                                               NaN
                                                                      NaN
                                                                             NaN
                                                                                    NaN ...
               water
                                  bar
                                       wheat rice
                                                    tea
          5 rows × 27 columns
```

```
In [119]: #as we can see from the result the 20th to 27th column has all null values
           store data.drop(store data.columns[[19,20,21,22,23,24,25,26]], axis=1, inpl
In [120]: store_data.isnull().sum()
Out[120]: 0
                     0
           1
                  1754
           2
                  3112
           3
                  4156
           4
                  4972
           5
                  5637
           6
                  6132
           7
                  6520
           8
                  6847
           9
                  7106
           10
                  7245
           11
                  7347
           12
                  7414
           13
                  7454
           14
                  7476
           15
                  7493
           16
                  7497
           17
                  7497
           18
                  7498
           dtype: int64
In [121]:
          store_data.head()
Out[121]:
                   0
                            1
                                     2
                                              3
                                                     4
                                                            5
                                                                 6
                                                                        7
                                                                               8
                                                                                      9
                                                                                            10
                                                        whole
                                                                                           low
                                                  green
                                                                           energy
                                       vegetables
                                                                    cottage
                                                                                  tomato
               shrimp
                       almonds avocado
                                                         weat yams
                                                                                            fat
            0
                                             mix grapes
                                                                    cheese
                                                                             drink
                                                                                   juice
                                                         flour
                                                                                         yogurt
               burgers meatballs
                                  eggs
                                            NaN
                                                   NaN
                                                         NaN
                                                               NaN
                                                                      NaN
                                                                             NaN
                                                                                    NaN
                                                                                          NaN
              chutney
                          NaN
                                  NaN
                                            NaN
                                                   NaN
                                                         NaN
                                                               NaN
                                                                      NaN
                                                                             NaN
                                                                                    NaN
                                                                                          NaN
            3
                       avocado
                                  NaN
                                            NaN
                                                   NaN
                                                         NaN
                                                               NaN
                                                                      NaN
                                                                             NaN
                                                                                    NaN
                                                                                          NaN
                turkey
                                           whole
               mineral
                                energy
                                                  green
                          milk
                                                         NaN
                                                               NaN
                                                                      NaN
                                                                             NaN
                                                                                    NaN
                                                                                          NaN
                water
                                   bar
                                        wheat rice
In [122]:
           #Data Proprocessing
           #The Apriori library we are going to use requires our dataset to be in the
           #where the whole dataset is a big list and each transaction in the dataset
           #Currently we have data in the form of a pandas dataframe.
           #To convert our pandas dataframe into a list of lists, we execute the scrip
           records = []
In [123]:
           for i in range(0, 7501):
```

records.append([str(store\_data.values[i,j]) for j in range(0, 18)])

```
In [124]: #Applying Apriori
          # The first parameter is the list of list that you want to extract rules fr
          # The second parameter is the min support parameter.
          # This parameter is used to select the items with support values greater the
          # Next, the min_confidence parameter filters those rules that have confiden
          # Similarly, the min_lift parameter specifies the minimum lift value for the
          # Finally, the min_length parameter specifies the minimum number of items t
In [125]: # Let's suppose that we want rules for only those items that are purchased
          # since our dataset is for a one-week time period.
          # The support for those items can be calculated as 35/7500 = 0.0045.
          # The minimum confidence for the rules is 20% or 0.2.
          # Similarly, we specify the value for lift as 3 and finally min length is 2
In [126]: | association_rules = apriori(records, min_support=0.0045, min_confidence=0.2
In [127]: # Viewing the Results
In [128]: | association results = list(association rules)
In [129]: |print(association_results)
          [RelationRecord(items=frozenset({'light cream', 'chicken'}), support=0.0
          04532728969470737, ordered statistics=[OrderedStatistic(items base=froze
          nset({'light cream'}), items_add=frozenset({'chicken'}), confidence=0.29
          059829059829057, lift=4.84395061728395)]), RelationRecord(items=frozense
          t({'escalope', 'mushroom cream sauce'}), support=0.005732568990801226, o
          rdered statistics=[OrderedStatistic(items base=frozenset({'mushroom crea
          m sauce'}), items_add=frozenset({'escalope'}), confidence=0.300699300699
          3007, lift=3.790832696715049)]), RelationRecord(items=frozenset({'past
          a', 'escalope'}), support=0.005865884548726837, ordered statistics=[Orde
          redStatistic(items_base=frozenset({'pasta'}), items_add=frozenset({'esca
          lope'}), confidence=0.3728813559322034, lift=4.700811850163794)]), Relat
          ionRecord(items=frozenset({'ground beef', 'herb & pepper'}), support=0.0
          15997866951073192, ordered_statistics=[OrderedStatistic(items_base=froze
          nset({'herb & pepper'}), items_add=frozenset({'ground beef'}), confidenc
          e=0.3234501347708895, lift=3.2919938411349285)]), RelationRecord(items=f
          rozenset({'ground beef', 'tomato sauce'}), support=0.005332622317024397,
          ordered_statistics=[OrderedStatistic(items_base=frozenset({'tomato sauc
          e'}), items_add=frozenset({'ground beef'}), confidence=0.377358490566037
          7, lift=3.840659481324083)]), RelationRecord(items=frozenset({'whole whe
In [130]: |len(association_results)
Out[130]: 45
In [131]: print(association_results[0])
          RelationRecord(items=frozenset({'light cream', 'chicken'}), support=0.0045
          32728969470737, ordered statistics=[OrderedStatistic(items base=frozenset
          ({'light cream'}), items_add=frozenset({'chicken'}), confidence=0.29059829
```

059829057, lift=4.84395061728395)])

```
In [132]: for item in association results:
            # first index of the inner list
            # Contains base item and add item
            pair = item[0]
            items = [x for x in pair]
            print("Rule: " + items[0] + " -> " + items[1])
            #second index of the inner list
            print("Support: " + str(item[1]))
            #third index of the list located at 0th
            #of the third index of the inner list
            print("Confidence: " + str(item[2][0][2]))
            print("Lift: " + str(item[2][0][3]))
            print("======="")
         Rule: light cream -> chicken
         Support: 0.004532728969470737
         Confidence: 0.29059829059829057
         Lift: 4.84395061728395
         _____
         Rule: escalope -> mushroom cream sauce
         Support: 0.005732568990801226
         Confidence: 0.3006993006993007
         Lift: 3.790832696715049
         _____
         Rule: pasta -> escalope
         Support: 0.005865884548726837
         Confidence: 0.3728813559322034
         Lift: 4.700811850163794
         _____
         Rule: ground beef -> herb & pepper
         Support: 0.015997866951073192
         Confidence: 0.3234501347708895
         Lift: 3.2919938411349285
```

In [ ]: