

DMW\_p6.ipynb

File Edit View Insert Runtime Tools Help All changes saved

Comment Share RAM Disk

## DMW Practical 6

### Question 1:

Assume a dataset of 10 transactions which has the list of the items that have been bought for each transaction. So the dataset will be having two kind of information in the dataset that is – (1) Transaction ID and (2) List of items. Find the support value of the each combination of the items.

```
[46] !pip install mlxtend
[47] import warnings
warnings.filterwarnings("ignore", category=DeprecationWarning)

[48] data1 = [['Milk', 'Onion', 'Nutmug', 'Kidney Beans', 'Eggs', 'Yogurt'],
           ['Milk', 'Onion', 'Nutmug', 'Kidney Beans', 'Eggs', 'Yogurt'],
           ['Milk', 'Apple', 'Kidney Beans', 'Eggs'],
           ['Milk', 'Unicorn', 'Corn', 'Kidney Beans', 'Yogurt'],
           ['Corn', 'Onion', 'Kidney Beans', 'Icecream', 'Eggs']]

[49] data1
[[{'item': 'Milk', 'itemset': 'Milk'}, {'item': 'Onion', 'itemset': 'Onion'}, {"item": "Nutmug", "itemset": "Nutmug"}, {"item": "Kidney Beans", "itemset": "Kidney Beans"}, {"item": "Eggs", "itemset": "Eggs"}, {"item": "Yogurt", "itemset": "Yogurt"}], [{"item": "Milk", "itemset": "Milk"}, {"item": "Onion", "itemset": "Onion"}, {"item": "Nutmug", "itemset": "Nutmug"}, {"item": "Kidney Beans", "itemset": "Kidney Beans"}, {"item": "Eggs", "itemset": "Eggs"}, {"item": "Yogurt", "itemset": "Yogurt"}], [{"item": "Milk", "itemset": "Milk"}, {"item": "Apple", "itemset": "Apple"}, {"item": "Kidney Beans", "itemset": "Kidney Beans"}, {"item": "Eggs", "itemset": "Eggs"}], [{"item": "Milk", "itemset": "Milk"}, {"item": "Unicorn", "itemset": "Unicorn"}, {"item": "Corn", "itemset": "Corn"}, {"item": "Kidney Beans", "itemset": "Kidney Beans"}, {"item": "Yogurt", "itemset": "Yogurt"}], [{"item": "Corn", "itemset": "Corn"}, {"item": "Onion", "itemset": "Onion"}, {"item": "Kidney Beans", "itemset": "Kidney Beans"}, {"item": "Icecream", "itemset": "Icecream"}, {"item": "Eggs", "itemset": "Eggs"}]]
```

```
[50] import pandas as pd
from mlxtend.preprocessing import TransactionEncoder

[51] encoder = TransactionEncoder()
encoded_data = encoder.fit(data1).transform(data1)
encoded_data
```

```
[52] df1 = pd.DataFrame(encoded_data,columns=encoder.columns_)
df1
```

	Apple	Corn	Dill	Eggs	Icecream	Kidney Beans	Milk	Nutmug	Onion	Unicorn	Yogurt
0	False	False	False	True	False	True	True	True	True	False	True
1	False	False	True	True	False	True	False	True	True	False	True
2	True	False	False	True	False	True	False	False	False	False	False
3	False	True	False	False	False	True	True	False	False	True	True
4	False	True	False	True	True	False	True	False	False	False	False

Next steps: [View recommended plots](#)

```
[53] test_data = encoded_data[:5]
encoder.inverse_transform(test_data)
encoded_data.astype('int')

array([[0, 0, 0, 1, 0, 1, 1, 1, 0, 1],
       [0, 0, 1, 1, 0, 1, 0, 1, 0, 1],
       [1, 0, 0, 1, 0, 1, 0, 0, 0, 0],
       [0, 1, 0, 0, 0, 1, 1, 0, 0, 1],
       [0, 1, 0, 1, 1, 0, 0, 1, 0, 1]])
```

```
[54] encoder.columns_
['Apple',
 'Corn',
 'Dill',
 'Eggs',
 'Icecream',
 'Kidney Beans',
 'Milk',
 'Nutmug',
 'Onion',
 'Unicorn',
 'Yogurt']
```

```
[55] from mlxtend.frequent_patterns import apriori
apriori(df1,min_support=0.4)
```

	support	itemsets
0	0.4	(1)
1	0.8	(3)
2	1.0	(5)
3	0.6	(6)
4	0.4	(7)
5	0.6	(8)
6	0.6	(10)
7	0.4	(1, 5)
8	0.8	(3, 5)
9	0.4	(3, 6)
10	0.4	(3, 7)
11	0.6	(8, 3)
12	0.4	(10, 3)
13	0.6	(5, 6)
14	0.4	(5, 7)
15	0.6	(8, 5)
16	0.6	(10, 5)
17	0.4	(10, 6)
18	0.4	(8, 7)
19	0.4	(10, 7)
20	0.4	(8, 10)
21	0.4	(3, 5, 6)
22	0.4	(3, 5, 7)
23	0.6	(8, 3, 5)
24	0.4	(10, 3, 5)
25	0.4	(8, 3, 7)
26	0.4	(10, 3, 7)
27	0.4	(8, 10, 3)
28	0.4	(10, 5, 6)
29	0.4	(8, 5, 7)
30	0.4	(10, 5, 7)

```
[56] apriori(df1,min_support=0.1)
```

	support	itemsets
0	0.2	(0)
1	0.4	(1)
2	0.2	(2)
3	0.8	(3)
4	0.2	(4)
...	...	...
144	0.4	(3, 5, 7, 8, 10)
145	0.2	(3, 6, 7, 8, 10)
146	0.2	(5, 6, 7, 8, 10)
147	0.2	(2, 3, 5, 7, 8, 10)
148	0.2	(3, 5, 6, 7, 8, 10)

149 rows × 2 columns

```
[57] apriori(df1, min_support = 0.4, use_colnames=True)
      /   0.4 (Kidney Beans, Onion)
```

	support	itemsets
0	0.2	(Kidney Beans, Onion)
1	0.8	(Kidney Beans, Eggs)
2	0.4	(Milk, Eggs)
10	0.4	(Nutmug, Eggs)
11	0.6	(Onion, Eggs)
12	0.4	(Yogurt, Eggs)
13	0.6	(Kidney Beans, Milk)
14	0.4	(Kidney Beans, Nutmeg)
15	0.6	(Kidney Beans, Onion)
16	0.6	(Kidney Beans, Yogurt)
17	0.4	(Milk, Yogurt)
18	0.4	(Onion, Nutmeg)
19	0.4	(Onion, Yogurt)
20	0.4	(Kidney Beans, Milk, Eggs)
21	0.6	(Kidney Beans, Onion, Eggs)
24	0.4	(Kidney Beans, Yogurt, Eggs)
25	0.4	(Onion, Nutmeg, Eggs)
26	0.4	(Onion, Yogurt, Eggs)
27	0.4	(Kidney Beans, Milk, Yogurt)
29	0.4	(Kidney Beans, Onion, Nutmeg)
30	0.4	(Kidney Beans, Nutmeg, Yogurt)
31	0.4	(Kidney Beans, Onion, Yogurt)
32	0.4	(Onion, Yogurt, Nutmeg)
33	0.4	(Kidney Beans, Onion, Nutmeg, Eggs)
34	0.4	(Kidney Beans, Nutmeg, Yogurt, Eggs)
35	0.4	(Kidney Beans, Onion, Yogurt, Eggs)
36	0.4	(Onion, Yogurt, Eggs, Nutmeg)
37	0.4	(Kidney Beans, Onion, Yogurt, Nutmeg)
38	0.4	(Onion, Nutmeg, Eggs, Yogurt, Kidney Beans)

```
[58] frequent_items = apriori(df1, min_support=0.4, use_colnames=True)
frequent_items['length'] = frequent_items['itemsets'].apply(lambda x: len(x))
```

```
[59] frequent_items
```

	support	itemsets	length
0	0.4	(Corn)	1
1	0.8	(Eggs)	1
2	1.0	(Kidney Beans)	1
3	0.6	(Milk)	1
4	0.4	(Nutmug)	1
5	0.6	(Onion)	1
6	0.6	(Yogurt)	1
7	0.4	(Kidney Beans, Corn)	2
8	0.8	(Kidney Beans, Eggs)	2
9	0.4	(Milk, Eggs)	2
10	0.4	(Nutmug, Eggs)	2
11	0.6	(Onion, Eggs)	2
12	0.4	(Yogurt, Eggs)	2
13	0.6	(Kidney Beans, Milk)	2
14	0.4	(Kidney Beans, Nutmeg)	2
15	0.6	(Kidney Beans, Onion)	2
16	0.6	(Kidney Beans, Yogurt)	2
17	0.4	(Milk, Yogurt)	2
18	0.4	(Onion, Nutmeg)	2
19	0.4	(Onion, Yogurt)	2
20	0.4	(Kidney Beans, Milk, Eggs)	3
21	0.4	(Kidney Beans, Onion, Eggs)	3
24	0.4	(Kidney Beans, Yogurt, Eggs)	3
25	0.4	(Onion, Nutmeg, Eggs)	3
26	0.4	(Onion, Yogurt, Eggs)	3
27	0.4	(Kidney Beans, Milk, Yogurt)	3
29	0.4	(Kidney Beans, Onion, Nutmeg)	3
30	0.4	(Kidney Beans, Nutmeg, Yogurt)	3

Next steps: [View recommended plots](#)

```
[60] frequent_items[frequent_items['itemsets'].apply(lambda x: 'Milk' in x and 'Eggs' in x)]
```

	support	itemsets	length
9	0.4	(Milk, Eggs)	2
21	0.4	(Kidney Beans, Milk, Eggs)	3

```
[61] frequent_items[frequent_items['itemsets'] == {'Milk', 'Eggs'}]
```

	support	itemsets	length
9	0.4	(Milk, Eggs)	2

```
[62] frequent_items[(frequent_items['length'] == 2) & (frequent_items['support'] >= 0.6)]
```

	support	itemsets	length
8	0.8	(Kidney Beans, Eggs)	2
11	0.6	(Onion, Eggs)	2
13	0.6	(Kidney Beans, Milk)	2
15	0.6	(Kidney Beans, Onion)	2
16	0.6	(Kidney Beans, Yogurt)	2

```
[63] import io
import pandas as pd
```

```
[64] data2 = pd.read_excel('OnlineRetail.xlsx')
data2.head()
```

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
536365	85123A	WHITE HANGING HEART T-LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	17850.0	United Kingdom
536365	71053	WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	2010-12-01 08:26:00	2.75	17850.0	United Kingdom
536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
536365	84029E	RED WOOLLY HOTIE WHITE HEART	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom

149 rows × 8 columns

1. As a data analyst, analyze the data and try to find the correlation between the items.

```
[65] from mlxtend.frequent_patterns import apriori
from mlxtend.frequent_patterns import association_rules
```

```
[66] basket = (data2.groupby(['InvoiceNo', 'Description'])['Quantity']
             .sum().unstack().reset_index().fillna(0)
             .set_index('InvoiceNo'))
```

```
[67] def encode_units(x):
    if x <= 0:
        return 0
    if x >= 1:
        return 1
    return 1
```

```
[68] basket_sets = basket.applymap(encode_units)
```

```
[69] frequent_items = apriori(basket_sets, min_support=0.05, use_colnames=True)
```

```
[70] print(data2.dtypes)
```

```
[71] numeric_data = data2.select_dtypes(include=['int64', 'float64'])
item_correlation = numeric_data.corr()
```

```
[72] high_price_items = data2[data2['UnitPrice'] > data2['UnitPrice'].quantile(0.95)].apply(lambda x: len(x))
```

```
[73] frequent_items[frequent_items['itemsets'].apply(lambda x: 'Milk' in x and 'Eggs' in x)]
```

	support	itemsets	length
9	0.4	(Milk, Eggs)	2

```
[74] frequent_items[frequent_items['itemsets'] == {'Milk', 'Eggs'}]
```

	support	itemsets	length
9	0.4	(Milk, Eggs)	2

```
[75] frequent_items[(frequent_items['length'] == 2) & (frequent_items['support'] >= 0.6)]
```

	support	itemsets	length
8	0.8	(Kidney Beans, Eggs)	2
11	0.6	(Onion, Eggs)	2
13	0.6	(Kidney Beans, Milk)	2
15	0.6	(Kidney Beans, Onion)	2
16	0.6	(Kidney Beans, Yogurt)	2

Next steps: [View recommended plots](#)

2. Find the most popular items among the customers.

```
[76] popular_items = data2['Description'].value_counts().head(10)
```

```
[77] popular_items
```

Description	Count
WHITE HANGING HEART T-LIGHT HOLDER	2369
REGENCY CAKESTAND 3 TIER	2200
JUMBO BAG RED RETROSPOT	2159
PARTY BIRDS	1727
LUNCH BAG RETROSPOT	1538
ASSORTED COLOUR BIRD ORNAMENT	1501
PACK OF 72 RETROSPOT CAKE CASES	1501
PARTY BIRDS	1491
REGENCY CAKE STAND 3 TIER	1474
UNIVERSITY BAGS	1464

```
[78] item_combinations = data2.groupby(['InvoiceNo'])['Description'].apply(list)
item_combinations_counts = item_combinations.value_counts().head(10)
```

```
[79] item_combinations_counts
```

Description	Count
WHITE HANGING HEART T-LIGHT HOLDER	2369
POSTAGE	1163
DOTCOM POSTAGE	696
CREAM SWEETHEARTS	553
RED RETROSPOT CAKE STAND	542
Name: count, dtype: int64	542

3. Which items' combinations are bought on most frequent basis?

```
[80] high_price_items = data2[data2['UnitPrice'] > data2['UnitPrice'].quantile(0.95)].apply(lambda x: len(x))
```

```
[81] high_price_items
```

Description	Count
REGENCY CAKESTAND 3 TIER	2194
POSTAGE	1163
DOTCOM POSTAGE	696
CREAM SWEETHEARTS	553
RED RETROSPOT CAKE STAND	542
Name: count, dtype: int64	542

4. Mention the items and combinations of the items whose sale should be more focused?

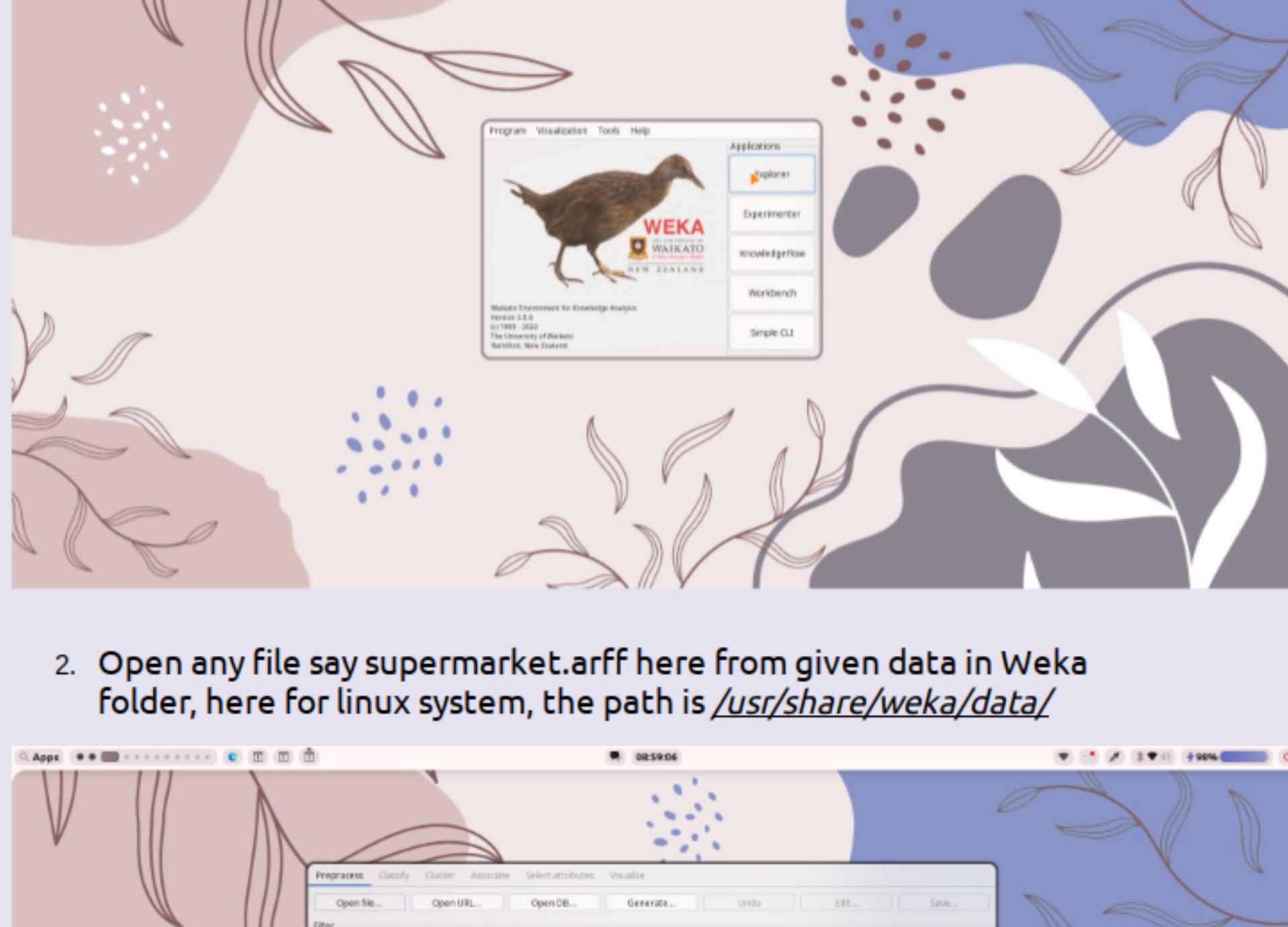
```
[82] high_price_items
```

Description	Count
REGENCY CAKESTAND 3 TIER	2194
POSTAGE	1163
DOTCOM POSTAGE	696
CREAM SWEETHEARTS	553
RED RETROSPOT CAKE STAND	542
Name: count, dtype: int64	542

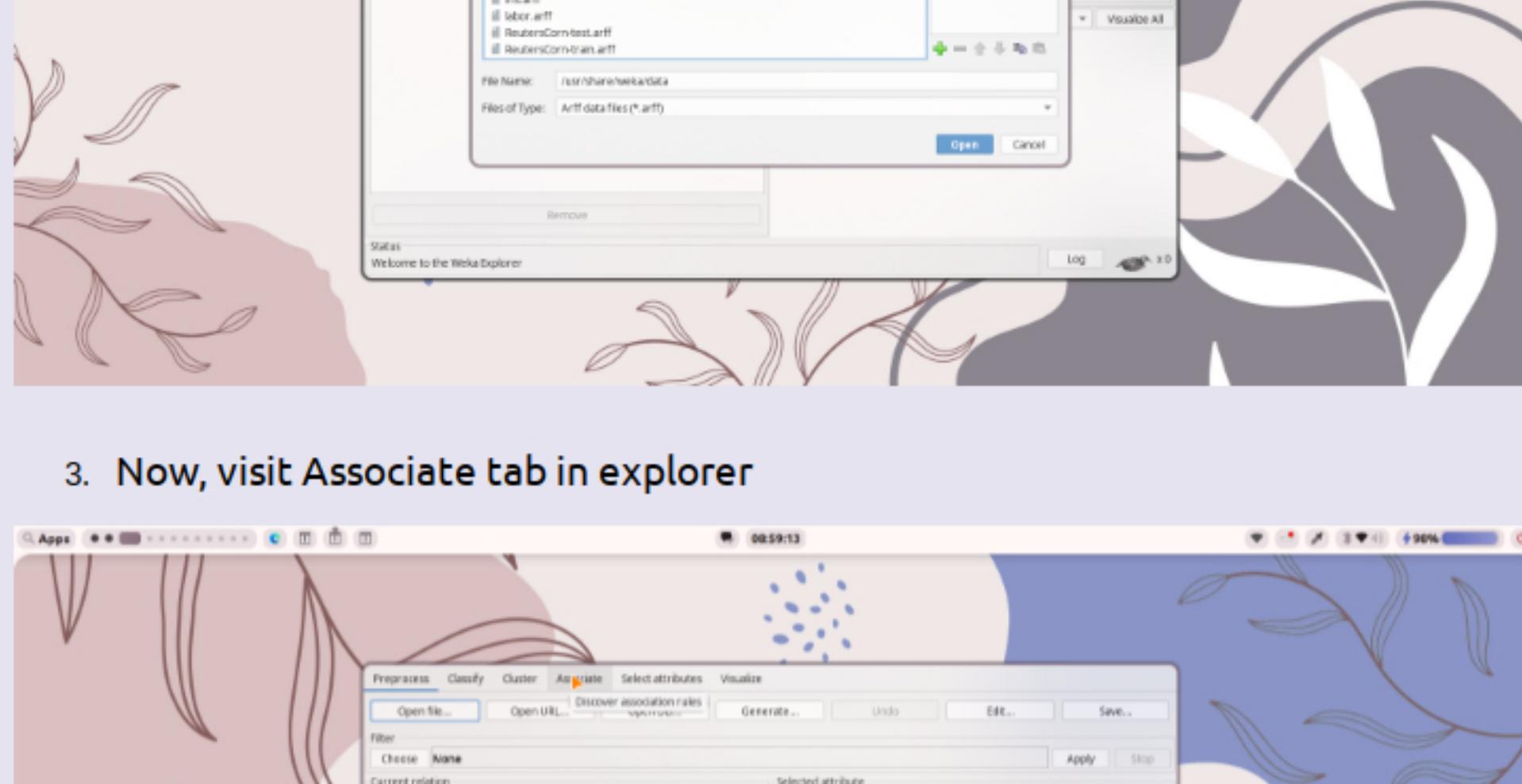
Name - Yash Lakhtariya  
Enrollment number - 21162101012  
Branch - CBA      Batch - 61  
DMW Practical 6

### Apriori testing using Weka tool :

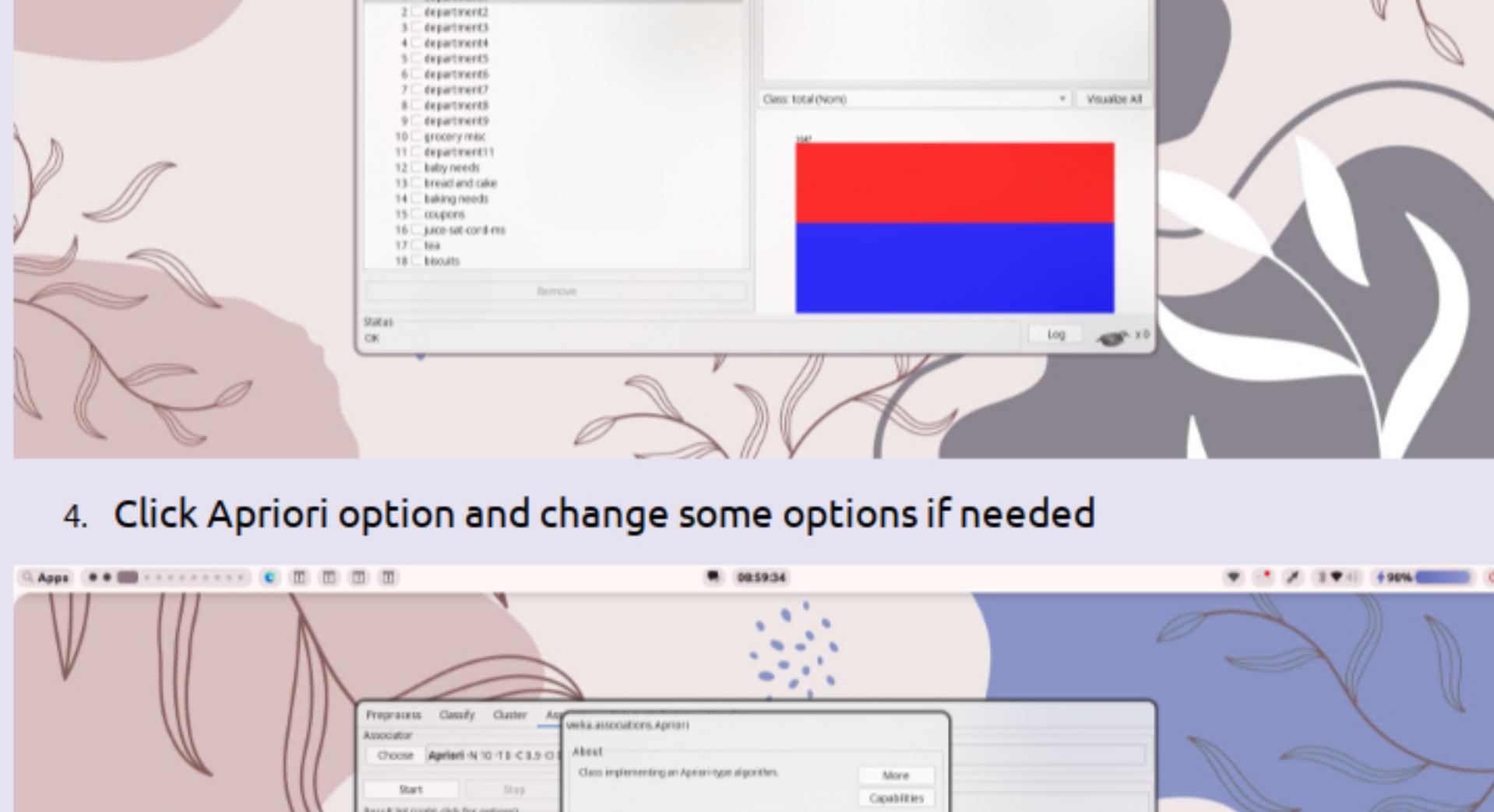
#### 1. Open Explorer in Weka tool



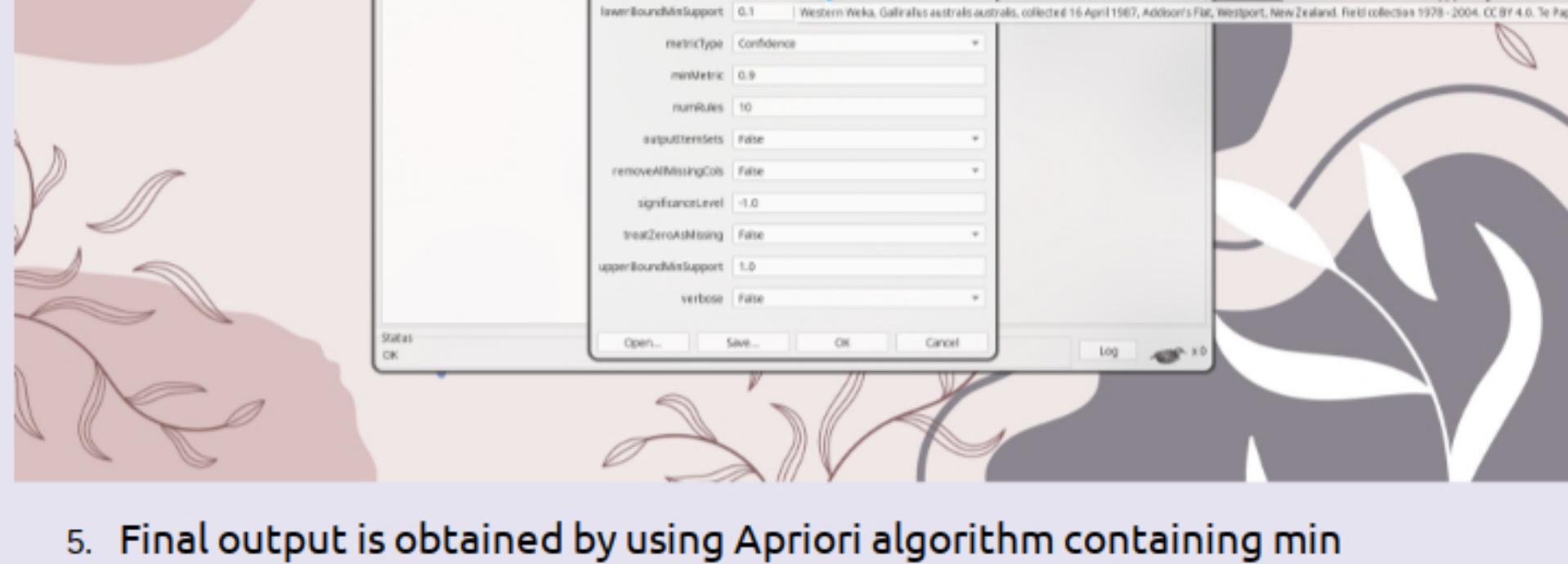
#### 2. Open any file say supermarket.arff here from given data in Weka folder, here for linux system, the path is /usr/share/weka/data/



#### 3. Now, visit Associate tab in explorer



#### 4. Click Apriori option and change some options if needed



#### 5. Final output is obtained by using Apriori algorithm containing min support, no. of cycles, sets of large items, best rules obtained

