

Worksheet__4b

PES University, Bangalore¶

Established under Karnataka Act No. 16 of 2013

UE20CS312 - Data Analytics - Worksheet 4b¶

Course instructor: Gowri Srinivasa, Professor Dept. of CSE, PES University

Designed by Harshith Mohan Kumar, Dept. of CSE - harshithmohankumar@pesu.pes.edu

Prerequisites¶

- Revise the following concepts
 - Boosting
 - * AdaBoost
 - Apriori Algorithm
- Install the following software
 - pandas
 - numpy
 - sklearn
 - matplotlib
 - mlxtend

Task¶

In this notebook you will be exploring how to implement and utilize AdaBoost and the Apriori algorithm. For AdaBoost, this notebook utilizes the standard dataset from sklearn. For Apriori, please ensure that you have downloaded the `BreadBasket_DMS.csv` within the same working directory.

Points¶

- Problem 1: 4 points
- Problem 2: 3 points
- Problem 3: 3 points

Loading the Dataset¶

In [1]:

```
# Imports
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
import pandas as pd
import numpy as np

# Load the wine dataset
data = datasets.load_wine(as_frame = True)

# Load x & y variables
X = data.data
y = data.target

# Split the dataset into train and test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 2)
```

Problem 1 (4 points)¶

Fit and evaluate the `AdaBoostClassifier` from `sklearn.ensemble` on the wine dataset. Use the `evaluate` model to print results.

Solution Steps:

1. From `sklearn.ensemble` import `AdaBoostClassifier`
2. Initialize the `AdaBoostClassifier` with `n_estimators` set to 30.
3. Use the `fit()` method and pass the train dataset.
4. Use the `evaluate(model, X_train, X_test, y_train, y_test)` method to print results.

For further reference: <https://www.kaggle.com/code/faressayah/ensemble-ml-algorithms-bagging-boosting-voting/notebook>

In [2]:

```
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report

# evaluate method to print results after training a particular model
def evaluate(model, X_train, X_test, y_train, y_test):
    y_test_pred = model.predict(X_test)
    y_train_pred = model.predict(X_train)

    print("TRAINING RESULTS: \n=====")
    clf_report = pd.DataFrame(classification_report(y_train, y_train_pred, output_dict=True))
```

```

print(f"CONFUSION MATRIX:\n{confusion_matrix(y_train, y_train_pred)}")
print(f"ACCURACY SCORE:\n{accuracy_score(y_train, y_train_pred):.4f}")
print(f"CLASSIFICATION REPORT:\n{clf_report}")

print("TESTING RESULTS: \n=====")
clf_report = pd.DataFrame(classification_report(y_test, y_test_pred, output_dict=True))
print(f"CONFUSION MATRIX:\n{confusion_matrix(y_test, y_test_pred)}")
print(f"ACCURACY SCORE:\n{accuracy_score(y_test, y_test_pred):.4f}")
print(f"CLASSIFICATION REPORT:\n{clf_report}")

```

Problem 2 (3 points)¶

Retrieve the frequent itemsets using the `apriori` method from `mlxtend.frequent_patterns`. The code below extracts the `basket_sets` and this is provided as input for the `apriori` method.

Solution Steps:

1. Use the `apriori` algorithm, set `min_support` to 0.03 and use `use_colnames` to `True`.
2. Print the output of the `apriori` method which provides the frequent itemsets

For further reference: <https://www.kaggle.com/code/victorcabral/bread-basket-analysis-apriori-association-rules/notebook> (Cells 26 onwards)

In [3]:

```

# Install mlxtend
!pip install mlxtend

```

```

Requirement already satisfied: mlxtend in /usr/local/lib/python3.8/dist-packages (0.21.0)
Requirement already satisfied: pandas>=0.24.2 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (1.3.4)
Requirement already satisfied: numpy>=1.16.2 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (1.21.0)
Requirement already satisfied: scipy>=1.2.1 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (1.7.3)
Requirement already satisfied: matplotlib>=3.0.0 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (3.5.1)
Requirement already satisfied: joblib>=0.13.2 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (1.1.0)
Requirement already satisfied: scikit-learn>=1.0.2 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (1.0.2)
Requirement already satisfied: setuptools in /usr/local/lib/python3.8/dist-packages (from mlxtend) (59.0.1)
Requirement already satisfied: python-dateutil>=2.8.1 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (2020.10.2)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (1.3.2)
Requirement already satisfied: pyparsing>=2.2.1 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (3.0.7)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (8.3.2)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (21.3)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (4.22.0)
Requirement already satisfied: cycycler>=0.10 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (0.10.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.8/dist-packages (from mlxtend) (2.2.0)
Requirement already satisfied: six>=1.5 in /usr/lib/python3/dist-packages (from python-dateutil) (1.16.0)

```

WARNING: You are using pip version 20.2.4; however, version 22.2.2 is available.
You should consider upgrading via the '/usr/bin/python3 -m pip install --upgrade pip' command

In [4]:

```
# Imports
from mlxtend.frequent_patterns import apriori
from mlxtend.frequent_patterns import association_rules
```

In [5]:

```
#Loading the dataset file
df = pd.read_csv('BreadBasket_DMS.csv')
```

In [6]:

```
df['Quantity'] = 1
df.head(7)
```

Out[6]:

	Date	Time	Transaction	Item	Quantity
0	2016-10-30	09:58:11	1	Bread	1
1	2016-10-30	10:05:34	2	Scandinavian	1
2	2016-10-30	10:05:34	2	Scandinavian	1
3	2016-10-30	10:07:57	3	Hot chocolate	1
4	2016-10-30	10:07:57	3	Jam	1
5	2016-10-30	10:07:57	3	Cookies	1
6	2016-10-30	10:08:41	4	Muffin	1

In [7]:

```
basket = df.groupby(['Transaction', 'Item'])['Quantity'].sum().unstack().fillna(0)
# There are a lot of zeros in the data but we also need to make sure any positive values are
# and anything less the 0 is set to 0. This step will complete the one hot encoding of the data
def encode_units(x):
    if x <= 0:
        return 0
    if x >= 1:
        return 1
```

```
basket_sets = basket.applymap(encode_units)
basket_sets
```

Out[7]:

Item

Adjustment

Afternoon with the baker

Alfajores

Argentina Night

Art Tray

Bacon

Baguette

Bakewell

Bare Popcorn

Basket

...

The BART

The Nomad

Tiffin

Toast

Truffles

Tshirt

Valentine's card

Vegan Feast

Vegan mincepie

Victorian Sponge

Transaction

1

0

0

0

0

0

0

0

0

0

0
0
0
0
...
0
0
0
0
0
0
0
0
0
0
0

9531 rows \times 95 columns

Problem 3 (3 points)¶

Now use the `association_rules` method and pass the `frequent_itemsets` as input (achieved using problem 2). Use `.head()` to display the top five rules.

Solution Steps:

1. Use the `association_rules` method, set metric to lift and `min_threshold` to 1.
2. Print the top five rules using `.head()`.

For further reference: <https://www.kaggle.com/code/victorcabral/bread-basket-analysis-apriori-association-rules/notebook> (Cell 32 and 33)