

PROJECT 1 - DATA SCIENCE

Authors:

Ariel Epshtein, ID: 316509504

Maor Mohav, ID: 316142363

Omer Ben David, ID: 316344449

Project structure:

- The project contains a dictionary that will save the values that will run the model using functions we have written, (in PYCHARM these values will be received as input from the user in the GUI, and in JUPYTER, the values are default, values we have chosen and can be changed)
- In each experiment, lead processing steps will be performed
- The project contains models and algorithms that use built-in libraries, but also models and algorithms with our implementation (functions with the extension at the end of the signature: `_Implementation`)
- Before running the project file in PYCHARM, the user will need to read the `Readme.txt` file
- We add documentation for all the functions in the Pytham project.

Libraries:

```
In [50]: import pandas as pd
import numpy as np
from sklearn.preprocessing import MinMaxScaler, StandardScaler
from scipy.stats import zscore
from sklearn.preprocessing import LabelEncoder
from pyutils import discrete, random, variable as drv
from sklearn import svm
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score
import cv2
import os
import random
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix, classification_report
from sklearn.cluster import KMeans
import pickle
from pandas_profiling import ProfileReport
```

Uploading data

```
In [78]: dic = {'train_path': 'D:\Vmedrive\OneDrive - ac.sce.ac.il\Desktop\train.csv', 'test_path': 'D:\Vmedrive\OneDrive - ac.sce.ac.il\Desktop\test.csv', 'missing_data': 'By all data', 'normalization': 'yes', 'discretization': '0'}
train_file = pd.read_csv(dic['train_path'])
test_file = pd.read_csv(dic['test_path'])
```

EDA:

```
In [79]: profile = ProfileReport(train_file, title='Pandas Profiling Report')
In [80]: profile
```

Pandas Profiling Report	Overview	Variables	Interactions	Correlations	Missing values	Sample
-------------------------	----------	-----------	--------------	--------------	----------------	--------

Overview

Dataset statistics	Variable types
Number of variables	10
Number of observations	42280
Missing cells	29
Missing cells (%)	<0.1%
Duplicate rows	0
Duplicate rows (%)	0.0%
Total size in memory	5.1 MB
Average record size in memory	120 B

Variables

age	Distinct	76	Minimum	18
Real number (R ₆₀)	Distinct (%)	0.2%	Maximum	500

Our analysis:

- Our data has 16 columns: 6 - numeric columns, 6 - categorical columns, 4 - boolean columns.
- Our data has 29 missing values.
- Columns correlations:
 - contact is highly correlated with month
 - education is highly correlated with job
 - job is highly correlated with education
 - month is highly correlated with contact and 2 other fields
 - day is highly correlated with month
 - month is highly correlated with housing
 - housing is highly correlated with month

Delete classification missing data

```
In [81]: data['class'].replace('', np.nan, inplace=True)
data.dropna(subset=['class'], inplace=True)
In [82]: Delete_Missing_Data
Delete_Missing_Data(test_file)
```

Complete Missing Data

```
In [83]: def numeric(i):
for i in i:
    if types[i] == str and i == '':
        return False
    return True
In [84]: def Q_Arrays(data):
a = []
no = data[data['class'] == 'no'].dropna(axis=0)
yes = data[data['class'] == 'yes'].dropna(axis=0)
for i in data.columns:
    if numeric(data[i]):
        a.append([yes[i].mean(), no[i].mean()])
    else:
        a.append([yes[i].value_counts().index[0], no[i].value_counts().index[0]])
return a
In [85]: def Missing_Data(train_file):
a = Q_Arrays(train_file)
c = data.columns.tolist()
if dic['missing_data'] == 'by all data':
    for i in data.columns:
        if numeric(data[i]) == True:
            data[i].fillna(train_file[i].mean(), inplace = True)
        else:
            data[i].fillna(train_file[i].value_counts().index[0], inplace = True)
    else:
        for i, rows in data.iterrows():
            for b in data.columns:
                if pd.isnull(data[b][i]) and data['class'][i] == 'yes':
                    data.at[i, str(b)] = a[c.index(b)][0]
                elif pd.isnull(data[b][i]) and data['class'][i] == 'no':
                    data.at[i, str(b)] = a[c.index(b)][1]
```

```
In [86]: Missing_Data(train_file)
Missing_Data(test_file)
```

Split the data

```
In [87]: for i in train_file.columns:
if numeric(train_file[i]):
    train_file[i] = pd.to_numeric(train_file[i], downcast = 'float')
    test_file[i] = pd.to_numeric(test_file[i], downcast = 'float')
```

Normalization

```
In [88]: def Normalization(train_file, test_file):
for b in train_file.columns:
    if dic['normalization'] == 'no':
        if numeric(train_file[b]):
            avg = train_file[b].mean()
            std = train_file[b].std()
            for i in range(0, len(train_file[i])):
                train_file.at[i, str(b)] = ((train_file[i][i] - avg) / std)
            for j in range(0, len(test_file[i])):
                test_file.at[j, str(b)] = ((test_file[i][i] - avg) / std)
In [89]: Normalization(train_file, test_file)
```

Discretization

```
In [90]: def Equal_Width_Binning(data, bins, labels):
for i in data:
    if numeric(data[i]):
        data[i] = pd.cut(data[i], bins, labels = label)
In [91]: def Equal_Frequency_Binning(data, bins):
for i in data:
    if numeric(data[i]):
        data[i] = pd.qcut(data[i], bins, duplicates = 'drop')
        best_feature = FeatureSelectByChi2(data[i].sort_values(), unique(), start=1)
        data[i] = data[i].cat.rename_categories(best_feature.categories)
In [92]: #Equal frequency discretization
def Equal_Frequency_Binning_Implementation():
def Equal_Frequency(data, test, columns, m):
    a = []
    for i in data.columns:
        a.append(i)
    a.sort()
    length = len(a)
    n = int(length / m)
    arr = []
    for i in range(0, m):
        for j in range(1, n, (i + 1) * n):
            if j == length:
                break
            arr = arr + [a[j]]
            arr.append(arr)
    for i in range(0, len(arr)):
        arr[i] = list(set(arr[i]))
    for i in range(0, len(test[columns])):
        for j in range(0, len(arr[i]):
            if test[columns][i] in arr[i]:
                test.at[i, str(columns)] = j
    for i in range(0, len(data[columns])):
        for j in range(0, len(arr[i]):
            if data[columns][i] in arr[i]:
                data.at[i, str(columns)] = j
    for i in train_file.columns:
        if numeric(train_file[i]):
            Equal_Frequency(train_file, test_file, i, dic['number_of_bins'])
def Equal_Width_Binning_Implementation():
def Equal_Width_Binning(data, test, columns, m):
    a = []
    for i in data.columns:
        a.append(i)
    a.sort()
    w = int(max(a) - min(a)) / m + 1
    min = min(a)
    arr = []
    for i in range(0, m + 1):
        arr = arr + [min + w * i]
        arr.append(arr)
    for i in range(0, m):
        temp = []
        for j in arr:
            if j == arr[i] and j <= arr[i+1]:
                temp = []
        arr.append(temp)
    for i in range(0, len(arr)):
        arr[i] = list(set(arr[i]))
    for i in range(0, len(test[columns])):
        for j in range(0, len(arr[i]):
            if test[columns][i] in arr[i]:
                test.at[i, str(columns)] = j
    for i in range(0, len(data[columns])):
        for j in range(0, len(arr[i]):
            if data[columns][i] in arr[i]:
                data.at[i, str(columns)] = j
    for i in train_file.columns:
        if numeric(train_file[i]):
            Equal_Width_Binning(train_file, test_file, i, dic['number_of_bins'])
In [93]: def Discretization(data):
bins = dic['number_of_bins']
label = []
for i in list(range(bins)):
    label.append(i)
if dic['discretization'] == 'Equal Width Binning':
    Equal_Width_Binning(data, bins, label)
elif dic['discretization'] == 'Equal Frequency Binning':
    Equal_Frequency_Binning(data, bins)
elif dic['discretization'] == 'Equal Width Binning-Implementation':
    Equal_Width_Binning_Implementation()
elif dic['discretization'] == 'Equal Frequency Binning-Implementation':
    Equal_Frequency_Binning_Implementation()
elif dic['discretization'] == 'Discretization based Entropy':
    for i in data.columns:
        if numeric(data[i]):
            new_val = np.array(data[i])
            H = -np.sum(new_val * np.log(new_val))
            for i in range(0, len(data[i])):
                if data[i][i] > 0:
                    else:
                        data.at[i, str(i)] = 1
                        data.at[i, str(i)] = 1
```

```
In [94]: Discretization(train_file)
Discretization(test_file)
```

Encoder

```
In [95]: def Encoder(data, dat):
for i in data.columns:
    if not numeric(data[i]):
        le = LabelEncoder()
        data[i] = le.fit_transform(data[i])
        dat[i] = le.fit_transform(dat[i])
In [96]: Encoder(train_file, test_file)
```

Saving files

```
In [97]: train_file.to_csv('train_file_clean.csv')
test_file.to_csv('test_file_clean.csv')
```

Model:

```
In [98]: def Naive_Bayes():
x_train = train_file.iloc[:, : -2]
y_train = train_file.iloc[:, -1]
x_test = test_file.iloc[:, : -1]
y_test = test_file.iloc[:, -1]
gnb = GaussianNB()
y_pred = gnb.fit(x_train, y_train).predict(x_test)
print('Number of misclassified points out of a total %d points : %d' % (X_test.shape[0], Y_test != y_pred).sum())
accuracy = accuracy_score(y_test, y_pred)*100
print('Precision = ', accuracy)
print('Confusion matrix:')
print(confusion_matrix(y_test, y_pred))
print('Report:')
print()
print(classification_report(y_test, y_pred))
pickle.dump(gnb, open('Naive_Bayes.sav', 'wb'))
In [99]: def Decision_Tree():
x_train = train_file.iloc[:, : -2]
y_train = train_file.iloc[:, -1]
x_test = test_file.iloc[:, : -1]
y_test = test_file.iloc[:, -1]
tree = DecisionTreeClassifier(criterion = 'entropy', random_state = 0)
tree.fit(x_train, y_train)
pred = tree.predict(x_test)
print('The prediction accuracy is: ', tree.score(x_test, y_test)*100, "%")
print('Confusion matrix:')
print(confusion_matrix(x_test, pred))
print('Report:')
print()
print(classification_report(y_test, pred))
print(classification_report(y_test, pred))
pickle.dump(tree, open('DT.sav', 'wb'))
In [100]: def Naive_Bayes_Implementation():
py = 0
pn = 0
total = len(test_file['class'])
correct = 0
y = 0
y = 0
for i in train_file['class']:
    if i == 0:
        y = 0 + 1
    else:
        y = y + 1
py = y / len(train_file['class'])
pn = n / len(train_file['class'])
for i, rows in test_file.iterrows():
    arryes = []
    arrno = []
    for j in test_file.columns:
        dyes = train_file[train_file[j] == test_file[i][i]]
        dno = train_file[train_file['class'] == 0]
        arryes.append(len(dyes['class']) / y)
        dno = train_file[train_file[j] == test_file[i][i]]
        dno = train_file[train_file['class'] == 0]
        arrno.append(len(dno['class']) / n)
    yes = 0
    no = 0
    for j in arryes:
        if yes < j:
            yes = j
        else:
            yes = yes * py
    for j in arrno:
        if no < j:
            no = j
        else:
            no = no * pn
    if yes > no and i == test_file['class'][i]:
        correct = correct + 1
    elif no > yes and i == test_file['class'][i]:
        correct = correct + 1
    print("The prediction accuracy is: %d" % (correct / total) * 100, "%")
    print("Accuracy of the model: ", (correct / len(test_file['class']))*100, "%")
```

```
In [101]: def Decision_Tree_Implementation():
dataset=train_file
test_dataset = test_file
def entropy(data_set):
    """this function calculate data set entropy"""
    Probability_set = []
    for i in data_set:
        counter = 0
        for j in data_set:
            if j == i:
                counter += 1
    Probability_set.append(counter / len(data_set))
    return sum(list(map(lambda x: -(x * log(x, 2)), Probability_set))) # return entropy
def InfoGain(data, split_attribute_name, target_name='class'):
    """InfoGain of the total dataset
    total_entropy = entropy(data[target_name])
    #Calculate the entropy of the dataset
    values, counts = np.unique(data[split_attribute_name], return_counts=True)
    W_Entropy = np.sum([counts[i]/np.sum(counts)*entropy(data.where(data[split_attribute_name]==values[i]).dropna()[target_name]) for i in range(len(values))])
    #Calculate the information gain
    Info_Gain = total_entropy - W_Entropy
    return Info_Gain
def ID3(data, originaldata, features, target_attribute='class', Parent_Node = None):
    if len(np.unique(data[target_attribute]) <= 1):
        return np.unique(data[target_attribute])[0]
    elif len(data) == 0:
        return np.unique(originaldata[target_attribute])[0]
    elif len(features) == 0:
        return Parent_Node
    else:
        Parent_Node = np.unique(data[target_attribute])[0].argmax(np.unique(data[target_attribute], return_counts=True)[1])
        item = [InfoGain(data, feature, target_attribute) for feature in features]
        index_of_BestFeature = np.argmax(item)
        BestFeature = features[index_of_BestFeature]
        #Create the tree structure. The root gets the name of the feature (best_feature) with the maximum information
        again in the first run
        tree = CreateFeatureTree()
        features = [i for i in features if i != BestFeature]
        for value in np.unique(data[BestFeature]):
            value = value
            sub_data = data where (data[BestFeature] == value) dropna()
            #Call the ID3 algorithm for each of those sub-records with the new parameters, here the recursion goes on!
            subtree = ID3(sub_data, subset, features, target_attribute, Parent_Node)
            tree[BestFeature][value] = subtree
        return (tree)
def predict(query, tree, default = 1):
    for key in list(query.keys()):
        if key in list(tree.keys()):
            try:
                result = tree[key][query[key]]
            except:
                return default
            result = tree[key][query[key]]
            if isinstance(result, dict):
                return predict(query, result)
            else:
                return result
def testing(data, tree):
    queries = data.iloc[:, :-2].to_dict(orient = "records")
    predicted = pd.DataFrame(columns=["predicted"])
    for i in range(len(data)):
        predicted.loc[i, 'predicted'] = predict(queries[i], tree, 1.0)
    print('The prediction accuracy is: ', (np.sum(predicted['predicted'] == data['class'])/len(data))*100, "%")
```

```
train_dataset=dataset.iloc[:100].reset_index(drop=True)
"""In the tree, Print the tree and predict the accuracy"""
tree = ID3(train_dataset, train_dataset, train_dataset.columns[:-1])
testing(train_dataset, tree)
In [102]: def Model(train_file, test_file):
if dic['model_type'] == 'Naive Bayes':
    Naive_Bayes()
elif dic['model_type'] == 'Naive Bayes-Implementation':
    Naive_Bayes_Implementation()
elif dic['model_type'] == 'Decision tree':
    Decision_Tree()
elif dic['model_type'] == 'Decision tree-Implementation':
    Decision_Tree_Implementation()
In [103]: Model(train_file, test_file)
Number of mislabeled points out of a total 3031 points : 1332
precision: 0.648679902029
confusion matrix:
[[129 180]
 [144 428]]
Report:
precision recall f1-score support
0 0.53 0.87 0.66 1467
1 0.69 0.27 0.39 1564
accuracy 0.61 0.57 0.62 3031
macro avg 0.61 0.57 0.62 3031
weighted avg 0.61 0.56 0.62 3031
```

KNN

```
In [104]: def Knn():
x_train = train_file.iloc[:, : -2]
y_train = train_file.iloc[:, -1]
x_test = test_file.iloc[:, : -1]
y_test = test_file.iloc[:, -1]
arr = []
for i in range(0, 10):
    classifier = KNeighborsClassifier(n_neighbors = i, p = 2, metric = 'euclidean')
    classifier.fit(x_train, y_train)
    y_pred = classifier.predict(x_test)
    err.append(accuracy_score(y_test, y_pred))
index = 0
ex = 0
for i in range(len(arr)):
    if arr[i] == ex:
        ex = arr[i]
        index = i
best_n_neighbors = (index + 2)
classifier = KNeighborsClassifier(n_neighbors = best_n_neighbors - 1, p = 2, metric = 'euclidean')
classifier.fit(x_train, y_train)
y_pred = classifier.predict(x_test)
print('The best n_neighbors is: ', best_n_neighbors)
print('The prediction accuracy is: ', ex * 100)
print('Confusion matrix:')
print(confusion_matrix(y_test, y_pred))
print('Report:')
print()
print(classification_report(y_test, y_pred))
pickle.dump(classifier, open('Knn.sav', 'wb'))
In [105]: Knn()
The best n_neighbors is: 3
The prediction accuracy is: 61.36229026592926
Confusion matrix:
[[142 165]
 [142 112]]
Report:
precision recall f1-score support
0 0.49 0.96 0.65 1467
1 0.63 0.67 0.65 1564
accuracy 0.56 0.53 0.59 3031
macro avg 0.56 0.53 0.59 3031
weighted avg 0.56 0.53 0.59 3031
```

K- means

```
In [106]: def K_means():
x_train = train_file.iloc[:, : -2]
y_train = train_file.iloc[:, -1]
x_test = test_file.iloc[:, : -1]
y_test = test_file.iloc[:, -1]
arr = []
for i in range(0, 10):
    classifier = KMeans(n_clusters = i)
    classifier.fit(x_train, y_train)
    y_pred = classifier.predict(x_test)
    err.append(accuracy_score(y_test, y_pred))
index = 0
ex = 0
for i in range(len(arr)):
    if arr[i] == ex:
        ex = arr[i]
        index = i
best_n_clusters = (index + 2)
classifier = KMeans(n_clusters = best_n_clusters)
classifier.fit(x_train, y_train)
y_pred = classifier.predict(x_test)
print('The best n_clusters is: ', best_n_clusters)
print('The prediction accuracy is: ', ex * 100)
print('Confusion matrix:')
print(confusion_matrix(y_test, y_pred))
print('Report:')
print()
print(classification_report(y_test, y_pred))
pickle.dump(classifier, open('K_means.sav', 'wb'))
In [107]: K_means()
The best n_clusters is: 3
The prediction accuracy is: 62.2829726162825
Confusion matrix:
[[145 162]
 [145 113]]
Report:
precision recall f1-score support
0 0.51 0.42 0.46 1467
1 0.63 0.62 0.57 1564
accuracy 0.52 0.52 0.52 3031
macro avg 0.52 0.52 0.52 3031
weighted avg 0.52 0.52 0.52 3031
```