# AI Computer Assignment #4

# Hossein Entezari Zarch

# 810196419

## Spring 2020

## ▾ Import Packages

At first we import needed packages on different parts of the assignment.

```
 1 import matplotlib.pyplot as plt
 2 import pandas as pd
 3 from sklearn.feature_selection import mutual_info_classif
 4 import numpy as np
 5 from sklearn.model_selection import train_test_split
 6 from sklearn.preprocessing import StandardScaler
 7 from sklearn.neighbors import KNeighborsClassifier
 8 from sklearn.metrics import classification_report, confusion_matrix
 9 from sklearn.tree import DecisionTreeClassifier
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
12 from sklearn.metrics import f1_score
13 from sklearn.ensemble import BaggingClassifier, RandomForestClassifier, VotingClassifier
```

## ▾ Read In the Data

In this part we read the data from .csv file with read_csv() function in pandas library, accoring the configuration present in parameters,

```
1 data = pd.read_csv("data.csv", true_values=["Yes"], false_values=["No"], index_col="Unnamed: 0")
```

## Phase0

## Data PreProcess

Saved successfully! ✕ taining information about country we defined two functions following two approaches:

-1. label: Through this approach we allocate a unique integer for each country and store it in the country column.

-2. one hot: Through this approach, we defined new columns for each counetry value we have and store the counrty value in a one hot and absolutely we omitted the counrty column finally.

Finally We decided to use one-hot method encoding in LogisticRegression & KNN alogorithms as they use the numeric value of feature arithmatic proccedure such as distance computaion in KNN and multiplying by weights in LogisticRegression, and we used the other m labeling in decision tree which totally based on that if the feature values are the same or not.

Date:In order to convert the date to processable values we convert the containing string to date and then with extracting the year and n and day values of the date we put them as three new attributes in data in order to get more information from date in our machine and p the overfitting.

Normalization: We defined normalize_data() function which performs the normalization process for us and it would help us in seeing a progressed classification algorithm.

At the end we defined a function named "preprocess_data" which calls the defined functions serially on the given data and puts the preprocessed data on the output.

Quantization: Also, as decision tree is based on labeled values, we decided to use labels for each period of data in 'Total Price' field wh values are fractional and we used quantization method on this field in decision tree to prevent overfitting.

```
1 def label_countries(data):
2     country_transform = data["Country"].value_counts().to_dict()
3     for i, country in enumerate(country_transform):
4         country_transform[country] = i+1
```

```python
 5     data.replace({"Country": country_transform}, inplace=True)
 6     return data
 7
 8 def quantize_total_prices(data):
 9     price_transform = data['Total Price'].value_counts().to_dict()
10     for i, period in enumerate(price_transform):
11         price_transform[period] = i+1
12     data.replace({"Total Price": price_transform}, inplace=True)
13     return data
```

Saved successfully!                           ×

```python
                           _dummies(data.Country))
17     enc = OneHotEncoder(handle_unknown='ignore')
18     enc_df = pd.DataFrame(enc.fit_transform(data[['Country']]).toarray())
19     data = data.join(enc_df)
20     data = data.drop('Country', axis=1)
21     return data
22
23 def seperate_date(data):
24     data['year'] = pd.to_datetime(data['Date']).dt.year
25     data['month'] = pd.to_datetime(data['Date']).dt.month
26     data['day'] = pd.to_datetime(data['Date']).dt.day
27     data = data.drop('Date', axis=1)
28     return data
29
30 def normalize_data(X):
31     x_columns = X.columns
32     scaler = StandardScaler()
33     scaler.fit(X)
34     X = scaler.transform(X)
35     X = pd.DataFrame(X, columns=x_columns)
36     return X
37
38 def preprocess_data(data, countries_label):
39     if countries_label == "label":
40         data = label_countries(data)
41         data = quantize_total_prices(data)
42     else:
43         data = one_hot_countries(data)
```

```
44
45      data = seperate_date(data)
46      y = data["Is Back"]
47      X = data.drop("Is Back", axis=1)
48      X = normalize_data(X)
49      return X, y
```

## ▾ Information Gain

Saved successfully!                    ✕      nding that how much the label is dependant on a specific input feature of data.

We defined the "information_gain(X,y)" function gives the data and its labels as input and gives us the information gain value of each da
feature with respect to the label of the data.

We plotted the results, on the first plot you see the values of information according to each of features which gives us a vision on how
important to label data and we see that the value according to the month is the highest and it means that, the label values are the most
dependant on month value of transaction among others.

The second plot shows us the information gains on different countries when we used the one hot method to encode country values. If
country has a high value of gain, it means that happing the transaction in that country or not gives us a high level of information on pre
the labels among other countries to be happened at.

```
 1 def information_gain(X, y):
 2      ans = [[], []]
 3      for col in X.columns:
 4          col_gain = mutual_info_classif(np.expand_dims(X[col].to_numpy(), 1), y)
 5          # ans.append([col, col_gain])
 6          ans[0].append(col)
 7          ans[1].append(col_gain)
 8      return ans
 9
10 def plot_compare(gains, labels = []):
11      plt.figure(figsize=(40, 10))
12      plt.rcParams.update({'font.size': 22})
13
14      for i, gain in enumerate(gains):
```

```
15          if i < len(labels):
16              plt.plot(gain[0], gain[1], label = labels[i])
17          else:
18              plt.plot(gain[0], gain[1])
19      plt.legend()
20      plt.grid(True)
21      plt.show()
22
23 def plot_lists(lists, labels=[], title='', xLabel='', yLabel=''):
24      plt.figure(figsize=(40, 10))
```

```
                              t.size': 22})
27      plt.rcParams['figure.figsize'] = [40, 10]
28
29      for i, li in enumerate(lists):
30          if i < len(labels):
31              plt.plot(range(len(li)), li, label=labels[i])
32          else:
33              plt.plot(range(len(li)), li)
34      # fig.suptitle(title, fontsize=24)
35      plt.xlabel(xLabel, fontsize=22)
36      plt.ylabel(yLabel, fontsize=22)
37      plt.legend()
38      plt.grid(True)
39      plt.show()
40
41 X, y = preprocess_data(data, "label")
42 gains = information_gain(X, y)
43 X2, y2 = preprocess_data(data, "one_hot")
44 gains2 = information_gain(X2, y2)
45
46 gains2_part1 = [gains2[0][:7], gains2[1][:7]]
47 gains2_part2 = [gains2[0][7:], gains2[1][7:]]
48 plot_compare([gains, gains2_part1], ["label", "one hot"])
49 plot_compare([gains2_part2])
```
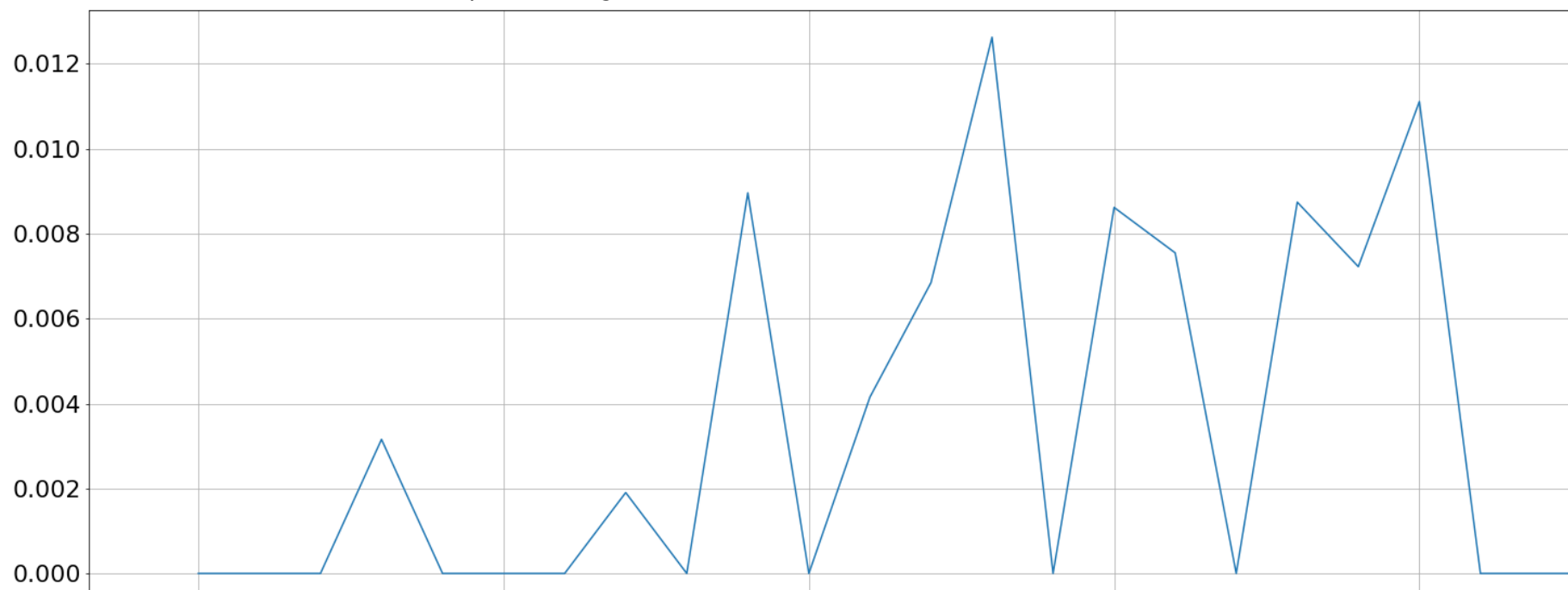
⟶

Saved successfully!

No handles with labels found to put in legend.

# Phase1

## Dataset Split

In this part with defining function "train_test_slit()" located in sklearn library, we split the dataset in to two parts containing 80 and 20 pe of original dataset.

> Saved successfully!    ✕   est = train_test_split(X, y, test_size=0.2)
>                             y2_test = train_test_split(X2, y2, test_size=0.2)

# Metrics Calculation

In this block we defined a function that gives us recall values for target labels and also precsion value and accuracy, these metrcis are computed with numbers stored in confusion matrix.

```
 1 def get_rpa_from_confusion(conf):
 2     num_true = conf[0][0] + conf[1][0]
 3     num_false = conf[1][0] + conf[1][1]
 4     rec_true = conf[0][0]/(conf[0][0]+conf[0][1])
 5     rec_false = conf[1][1]/(conf[1][1]+conf[1][0])
 6     recall = ((rec_true * num_true) + (rec_false * num_false))/(num_true + num_false)
 7     prec_true = conf[0][0]/(conf[0][0]+conf[1][0])
 8     prec_false = conf[1][1]/(conf[1][1]+conf[0][1])
 9     precision = ((prec_true * num_true) + (prec_false * num_false))/(num_true + num_false)
10     acc = (conf[0][0]+conf[1][1])/(conf[0][0]+conf[0][1]+conf[1][0]+conf[1][1])
11     return [recall, precision, acc]
```

# KNN

Metrics Calculation: We calculated recall and precision values on both train and test datasets by confusion matrix and plotted all of the we calcualted the accuracy and F1-Score on both train and test dataset and we found it useful in analyze the high bias and high varianc

Plot1: The values of precision and recall values on target labels on train and test dataset are plotted.

Plot2: The values of accuracy on both train and test dataset are plotted.

Plot3: The values of F1-Score on both train and test dataset are plotted.

Choose The Best HyperParameter: In order to choose the best num of neighbors for our model we see the two last plots wich plot the a and F1-Score of the model on both train & test dataset for each value of n-neighbors, we see that for n-neighbors above 20 we do not s dramatic dhange in accuracy and until 20 the train and test accuracies tend to about a same value near 72% so we choose [n-neighbor

Saved successfully! × ter. with little values of n-neighbor we see a high value of accuracy on train data and a low level c us that the model is overfitted.

Recall & Precision: We see that the recall and precision values are about the same values on this configuration and it means that our m not overfitted and with habing a high accuracy it is not also high bias, so it low bias and low variance.

```
1 recalls = [[], []]
2 precisions = [[], []]
3 accuracies = [[], []]
4 f1 = [[], []]
5
6 for i in range(1, 50):
7     classifier_KNN = KNeighborsClassifier(n_neighbors=i)
8     classifier_KNN.fit(X2_train, y2_train)
9
10    y2_train_pred = classifier_KNN.predict(X2_train)
11    conf_matrix = confusion_matrix(y2_train, y2_train_pred)
12    f1[0].append(f1_score(y2_train, y2_train_pred))
13
14    if i == 1:
15        print('total label nums', conf_matrix[0][0]+conf_matrix[0][1], conf_matrix[1][0]+conf_matrix[1][1])
16    rpa = get_rpa_from_confusion(conf_matrix)
17    recalls[0].append(rpa[0])
18    precisions[0].append(rpa[1])
19    accuracies[0].append(rpa[2])
20
21    y2_test_pred = classifier_KNN.predict(X2_test)
22    conf_matrix = confusion_matrix(y2_test, y2_test_pred)
```

```
23     f1[1].append(f1_score(y2_test, y2_test_pred))

24

25     if i == 1:
26         print('total label nums', conf_matrix[0][0]+conf_matrix[0][1], conf_matrix[1][0]+conf_matrix[1][1])
27     rpa = get_rpa_from_confusion(conf_matrix)
28     recalls[1].append(rpa[0])
29     precisions[1].append(rpa[1])
30     accuracies[1].append(rpa[2])

31

ls[1], precisions[0], precisions[1]], labels=['recalls_tain', 'recalls_test', 'pr
```

Saved successfully! ✕

```
35 plot_lists(accuracies, labels=['train', 'test'], title='accuracy')
36 print('F1-score')
37 plot_lists(f1, labels=['train', 'test'], title='f1')
```
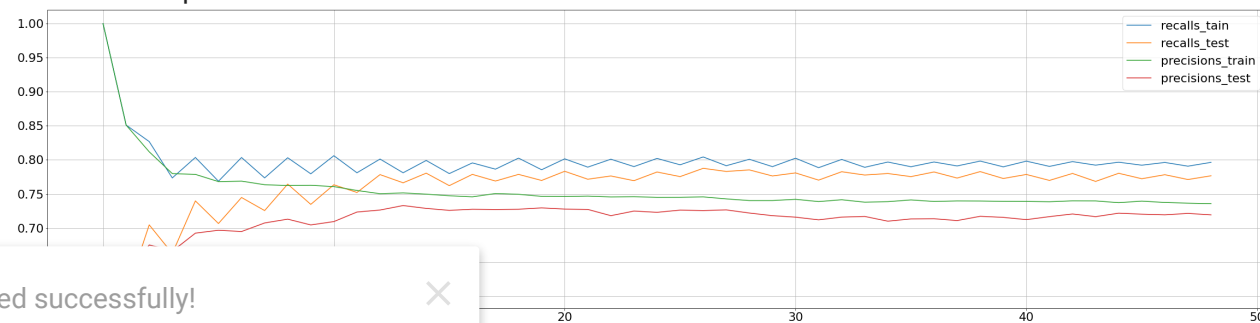
⤷

```
total label nums 1081 2411
total label nums 299 575
recalls & precisions
```
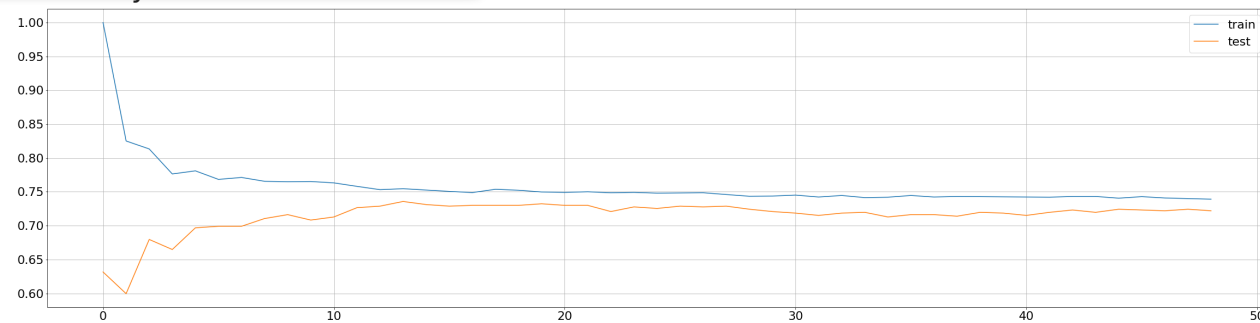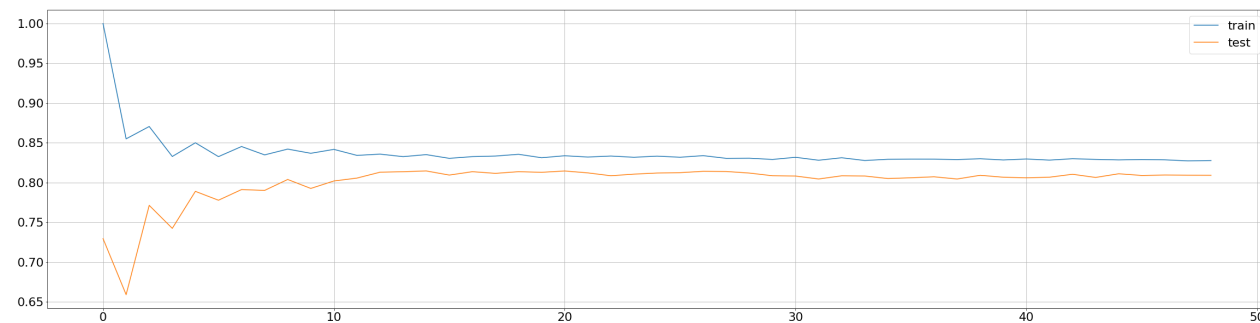


Saved successfully!                                    ×

Accuracy



F1-score

# ▾ Decision Tree

Choose the Best HyperParameter: The last plot shows us that, for max_depth of one the train and test accuracies are so near each othe are about 75% and after that the train starts to increase and the test starts to decrease and the model overfits, so we choose [max_dep

Saved successfully! ✕

recall & precision: with max_depth of one, the recall and precision values are not near each other and the model is biased, but with inc the max_depth they start to tend to one on train dataset and they tend to different values on test dataset.

```
 1 recalls = [[], []]
 2 precisions = [[], []]
 3 accuracies = [[], []]
 4
 5 for max_dpth in range(1, 50):
 6     classifier_tree = DecisionTreeClassifier(criterion="entropy", max_depth=max_dpth)
 7     classifier_tree.fit(X_train, y_train)
 8
 9     y_train_pred = classifier_tree.predict(X_train)
10     rpa = get_rpa_from_confusion(confusion_matrix(y_train, y_train_pred))
11     recalls[0].append(rpa[0])
12     precisions[0].append(rpa[1])
13     accuracies[0].append(rpa[2])
14
15     y_test_pred = classifier_tree.predict(X_test)
16     rpa = get_rpa_from_confusion(confusion_matrix(y_test, y_test_pred))
17     recalls[1].append(rpa[0])
18     precisions[1].append(rpa[1])
19     accuracies[1].append(rpa[2])
20
21 print('recalls & precisions')
22 plot_lists([recalls[0], recalls[1], precisions[0], precisions[1]], labels=['recalls_tain', 'recalls_test', 'pr
23 print('Accuracy')
24 plot_lists(accuracies, labels=['train', 'test'], title='accuracy')
```

```
24 plot_lists(accuracies, labels=['train', 'test'], title='accuracy')
```
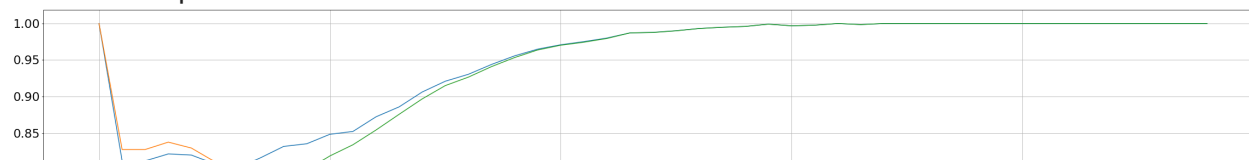
Saved successfully!                                            ✕

```
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:7: RuntimeWarning: invalid value encountered in
  import sys
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:7: RuntimeWarning: invalid value encountered in
  import sys
recalls & precisions
```



Saved successfully!                              ✕

In this model, using the normalized data is so important because it uses the values of features in multiplying them with the weights, an using the one-hot method to encode categorical feature plays a crucial role.

Accuracy: We see that the accuracy on both train and test dataset are both near 75%, and it is not absolutely in a overfit mode because also gives a high accuracy on train data.

Recall & Precision: Recall and Precision values are not completely near each other and it means that the model is biased on 'true' targe which is more popular to happen.



```
 1 classifierLogReg = LogisticRegression(random_state=1).fit(X2_train, y2_train)
 2
 3 y2_train_pred = classifierLogReg.predict(X2_train)
 4 rpa = get_rpa_from_confusion(confusion_matrix(y2_train, y2_train_pred))
 5 print('train--------------------------')
 6 print('recall train', rpa[0])
 7 print('precisions train', rpa[1])
 8 print('accuracy train', rpa[2])
 9
10 y2_test_pred = classifierLogReg.predict(X2_test)
11 rpa = get_rpa_from_confusion(confusion_matrix(y2_test, y2_test_pred))
12 print('\ntest-------------------------')
13 print('recall test', rpa[0])
14 print('precision test', rpa[1])
15 print('accuracy test', rpa[2])
```

```
train--------------------------
recall train 0.7777971675951405
precisions train 0.7435473203433672
accuracy train 0.7468499427262314

test--------------------------
recall test 0.7776457776457777
precision test 0.740334084084084
accuracy test 0.7425629290617849
```

Saved successfully! ✕

## ▾ Bagging KNN

n_estimators: We observed less values of n_estimators gives us better unbiased model with being the recall and precision values near other, finaaly we chose 5 as the value of this hyperparameter. also the mod

The accuracy metrics are given below the block.

```
 1 base_classifier = KNeighborsClassifier(n_neighbors=22)
 2 model = BaggingClassifier(base_estimator=base_classifier, n_estimators=5, max_samples=0.5, max_features=0.5, r
 3 model.fit(X2_train, y2_train)
 4
 5 y2_train_pred = model.predict(X2_train)
 6 y2_test_pred = model.predict(X2_test)
 7
 8 rpa = get_rpa_from_confusion(confusion_matrix(y2_train, y2_train_pred))
 9 print('train--------------------------')
10 print('recall train', rpa[0])
11 print('precisions train', rpa[1])
12 print('accuracy train', rpa[2])
13
14 rpa = get_rpa_from_confusion(confusion_matrix(y2_test, y2_test_pred))
15 print('\ntest--------------------------')
16 print('recall test', rpa[0])
17 print('precision test', rpa[1])
```

```
18 print('accuracy test', rpa[2])
19
```

```
train--------------------------
recall train 0.829301332548038
precisions train 0.73551863493489
accuracy train 0.7376861397479955

test--------------------------
recall test 0.8170759675107502
                              970818
                              91991
```

Saved successfully!

## Bagging Decision Tree:

The base model is the best model configured we could build in last phase and now we have 100 of them in a BaggingClassifier and we the accuracy is nearly the same as the accuracy of the base model, bagging algorithm could work on an overfit model and in a model th biscally overfitted we would not see a dramatic change or progress.

```
 1 base_classifier = DecisionTreeClassifier(criterion="entropy", max_depth=6)
 2 model = BaggingClassifier(base_estimator=base_classifier, n_estimators=100, max_samples=0.5, max_features=0.5,
 3 model.fit(X_train, y_train)
 4
 5 y_train_pred = model.predict(X_train)
 6 rpa = get_rpa_from_confusion(confusion_matrix(y_train, y_train_pred))
 7 print('train--------------------------')
 8 print('recall train', rpa[0])
 9 print('precisions train', rpa[1])
10 print('accuracy train', rpa[2])
11
12 y_test_pred = model.predict(X_test)
13 rpa = get_rpa_from_confusion(confusion_matrix(y_test, y_test_pred))
14 print('\ntest--------------------------')
15 print('recall test', rpa[0])
16 print('precision test', rpa[1])
17 print('accuracy test', rpa[2])
```

```
train--------------------------
recall train 0.8680152919027964
precisions train 0.7616874680951521
accuracy train 0.7597365406643757

test--------------------------
recall test 0.8496065470569719
precision test 0.7355604808525005
accuracy test 0.7368421052631579
```

Saved successfully! ✕

We focused on two hyperparameter n_estimators, max_depth and we plotted the results based on these hyperparameters and the metr
below. We observe from the last plot that, model with n_estimators=8 and max_depth=7, and with this configuration we have a model tl
overfit and its accuracy is about 76%.

```
1 recalls = [[], []]
2 precisions = [[], []]
3 accuracies = [[], []]
4
5 for n in range(1, 11):
6     model = RandomForestClassifier(n_estimators=n, criterion='entropy')
7     model.fit(X_train, y_train)
8
9     y_train_pred = model.predict(X_train)
10    rpa = get_rpa_from_confusion(confusion_matrix(y_train, y_train_pred))
11    recalls[0].append(rpa[0])
12    precisions[0].append(rpa[1])
13    accuracies[0].append(rpa[2])
14
15    y_test_pred = model.predict(X_test)
16    rpa = get_rpa_from_confusion(confusion_matrix(y_test, y_test_pred))
17    recalls[1].append(rpa[0])
18    precisions[1].append(rpa[1])
19    accuracies[1].append(rpa[2])
20
21 print('recalls & precisions')
```

```
22 plot_lists([recalls[0], recalls[1], precisions[0], precisions[1]], labels=['recalls_true', 'recalls_false', 'p
23 print('Accuracy')
24 plot_lists(accuracies, labels=['train', 'test'], title='accuracy')
25
26
27 recalls = [[], []]
28 precisions = [[], []]
29 accuracies = [[], []]
30
```

Saved successfully!　　✕

```
                        ifier(n_estimators=8, criterion='entropy', max_depth=m)
                        )
34
35     y_train_pred = model.predict(X_train)
36     rpa = get_rpa_from_confusion(confusion_matrix(y_train, y_train_pred))
37     recalls[0].append(rpa[0])
38     precisions[0].append(rpa[1])
39     accuracies[0].append(rpa[2])
40
41     y_test_pred = model.predict(X_test)
42     rpa = get_rpa_from_confusion(confusion_matrix(y_test, y_test_pred))
43     recalls[1].append(rpa[0])
44     precisions[1].append(rpa[1])
45     accuracies[1].append(rpa[2])
46
47 print('recalls & precisions')
48 plot_lists([recalls[0], recalls[1], precisions[0], precisions[1]], labels=['recalls_tain', 'recalls_test', 'pr
49 print('Accuracy')
50 plot_lists(accuracies, labels=['train', 'test'], title='accuracy')
```
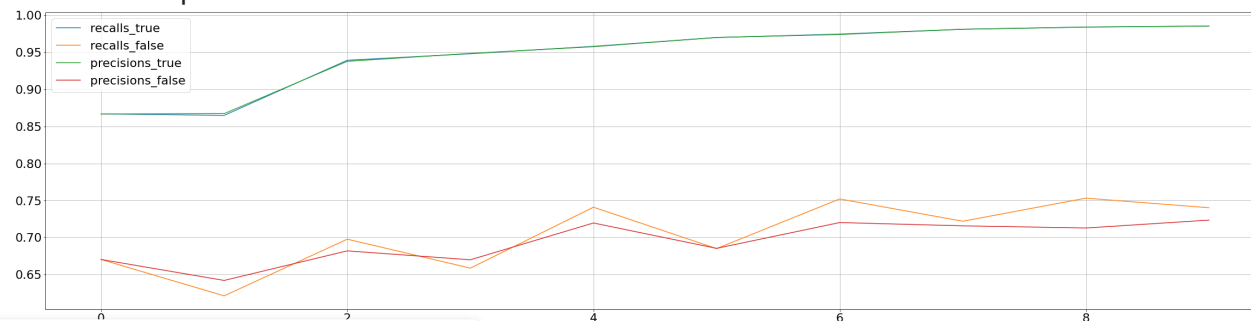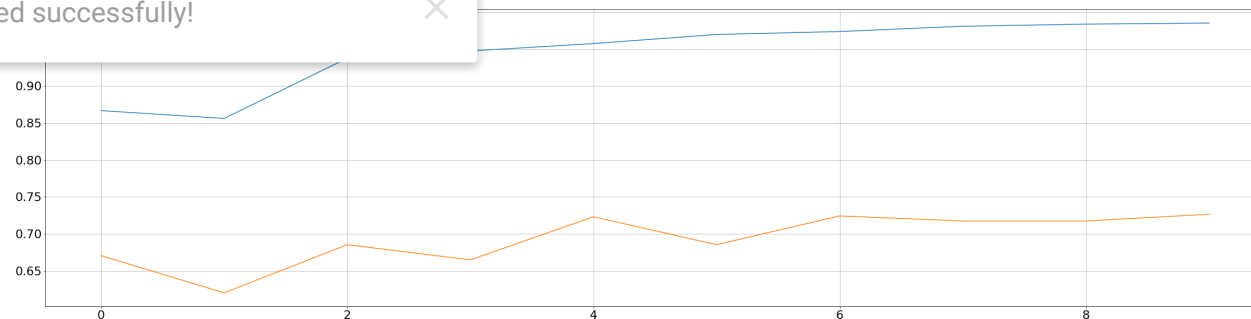
⊏→

## recalls & precisions



Saved successfully! ✕



```
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:7: RuntimeWarning: invalid value encountered in
  import sys
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:7: RuntimeWarning: invalid value encountered in
  import sys
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:7: RuntimeWarning: invalid value encountered in
  import sys
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:7: RuntimeWarning: invalid value encountered in
  import sys
```
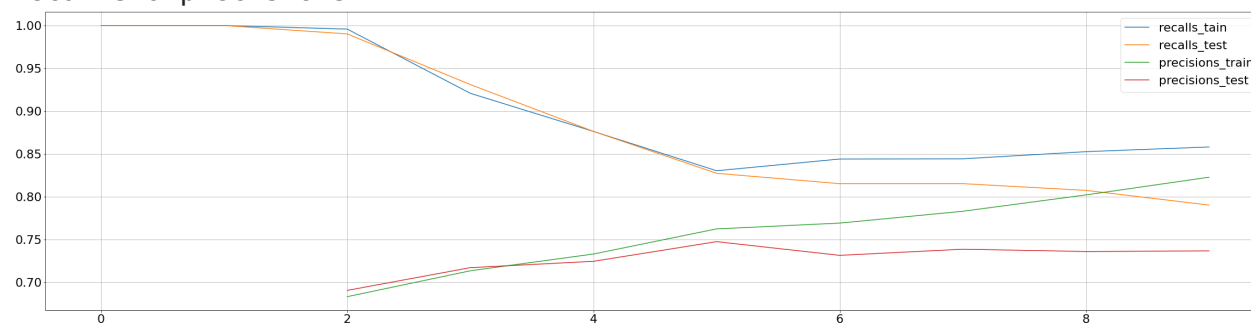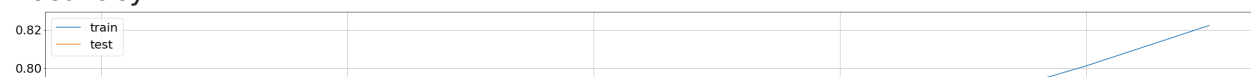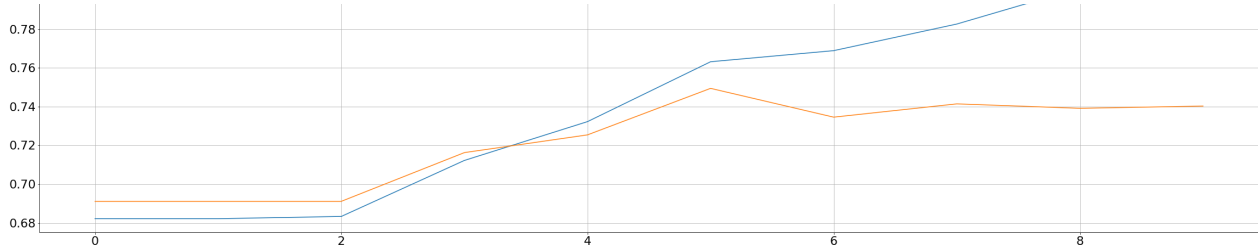
## recalls & precisions



## Accuracy

Saved successfully! ✕

## ▾ Solve Overfitting by Bagging

Overfitting: The base model is an overfit one as you can observe in its plot, but when we used the bagging model with n_estimators 1, w

_____ ted, and also when we used higher n_estimators we see that, both accuracies increased, but the

Saved successfully!　　　　　　　　　✕ so increased but, as finally with [n_estimators = 10] we had accuracy of 71% on test step and it w

higher than this value in n_estimators = 1, we decided this value to be 10 and also we think it is generally a benefcial vlue for this

hyperparameter.

```
 1 base_classifier = DecisionTreeClassifier(criterion="entropy", max_depth=15)
 2 model = BaggingClassifier(base_estimator=base_classifier, n_estimators=30, max_samples=0.5, max_features=0.5,
 3 model.fit(X_train, y_train)
 4
 5 y_train_pred = model.predict(X_train)
 6 rpa = get_rpa_from_confusion(confusion_matrix(y_train, y_train_pred))
 7 print('train-------------------------')
 8 print('recall train', rpa[0])
 9 print('precisions train', rpa[1])
10 print('accuracy train', rpa[2])
11
12 y_test_pred = model.predict(X_test)
13 rpa = get_rpa_from_confusion(confusion_matrix(y_test, y_test_pred))
14 print('\ntest-------------------------')
15 print('recall test', rpa[0])
16 print('precision test', rpa[1])
17 print('accuracy test', rpa[2])
```

⤷

```
train--------------------------
recall train 0.9149879774038698
precisions train 0.8905732501423107
accuracy train 0.8874570446735395
```

## BootStrapping

Through this approach, we sample subset of our data with replacement multiple times and after generating these subsets, we learn mu

mple datasets, this approach causes our model to overfit harder and we prevent it from overfittin

Saved successfully! ✕ a subset of the features of data on each model this approach can cause the bias to get higher.

## ▾ HardVoting:

In this approach, we defined our three classifier models and then put them in the hard-voting model and after fitting the model on the d classified the input data based on the class that is voted as the most among all of given classifiers.

```
 1 clf1 = KNeighborsClassifier(n_neighbors=22)
 2 clf2 = DecisionTreeClassifier(criterion="entropy", max_depth=6)
 3 clf3 = LogisticRegression(random_state=1)
 4
 5 model = VotingClassifier(estimators=[('knn', clf1), ('dt', clf2), ('lr', clf3)])
 6 model.fit(X2_train, y2_train)
 7
 8 y2_train_pred = model.predict(X2_train)
 9 rpa = get_rpa_from_confusion(confusion_matrix(y2_train, y2_train_pred))
10 print('train--------------------------')
11 print('recall train', rpa[0])
12 print('precisions train', rpa[1])
13 print('accuracy train', rpa[2])
14
15 y2_test_pred = model.predict(X2_test)
16 rpa = get_rpa_from_confusion(confusion_matrix(y2_test, y2_test_pred))
17 print('\ntest--------------------------')
18 print('recall test', rpa[0])
19 print('precision test', rpa[1])
```

```
20 print('accuracy test', rpa[2])
```

```
train--------------------------
recall train 0.7988795814444122
precisions train 0.7500179411475157
accuracy train 0.7528636884306987

test--------------------------
recall test 0.7964354867100863
precision test 0.7492972271025895
                               08627
```

Saved successfully!                          ✕

## Similarity of Models

In the result shown below you can see the ratio of the answers that are the to all of the answers, and as you see always it is 86% to 90% percent.

```
 1 clf1 = KNeighborsClassifier(n_neighbors=22)
 2 clf2 = DecisionTreeClassifier(criterion="entropy", max_depth=6)
 3 clf3 = LogisticRegression(random_state=1)
 4
 5 clf1.fit(X2_train, y2_train)
 6 clf2.fit(X2_train, y2_train)
 7 clf3.fit(X2_train, y2_train)
 8
 9 y2_train_pred1 = clf1.predict(X2_train)
10 y2_train_pred2 = clf2.predict(X2_train)
11 y2_train_pred3 = clf3.predict(X2_train)
12 diff_train12 = sum(y2_train_pred1 == y2_train_pred2)/len(y2_train_pred1)
13 diff_train13 = sum(y2_train_pred1 == y2_train_pred3)/len(y2_train_pred1)
14 diff_train23 = sum(y2_train_pred2 == y2_train_pred3)/len(y2_train_pred1)
15 print('train', '1-2', diff_train12, '1-3', diff_train13, '2-3', diff_train23)
16
17 y2_test_pred1 = clf1.predict(X2_test)
18 y2_test_pred2 = clf2.predict(X2_test)
19 y2_test_pred3 = clf3.predict(X2_test)
20 diff_test12 = sum(y2_test_pred1 == y2_test_pred2)/len(y2_test_pred1)
```

```
21 diff_test13 = sum(y2_test_pred1 == y2_test_pred3)/len(y2_test_pred1)
22 diff_test23 = sum(y2_test_pred2 == y2_test_pred3)/len(y2_test_pred1)
23 print('test', '1-2', diff_test12, '1-3', diff_test13, '2-3', diff_test23)
```

```
train 1-2 0.8791523482245132 1-3 0.9146620847651775 2-3 0.8883161512027491
test 1-2 0.8672768878718535 1-3 0.8913043478260869 2-3 0.8707093821510298
```

## Analyse Ensmble Methods:

Saved successfully!  ×  ke seeing less amount of data and being less complex prevent overfitting efficiently as you can s
~~~~~~~~~~~~~~~~~~~~~~~~~~~~gging" of the report, but in the case where the base models were not overfitted, the ensmble metho
not show any dramatic progress and it could not play a crucial role.

Also, in comparision of the prediction of models we observe that they about 90% on the same answer, so making vote between these n
can not cause a high level of progress in the model rather than base models theirselves.

```
1
```

Saved successfully!                                                    ✕