



MIDDLE EAST TECHNICAL UNIVERSITY
NORTHERN CYPRUS CAMPUS

CNG 562
MACHINE LEARNING

MIDTERM 1

Report

Kaan Taha Köken 2152064

April 21, 2020

Contents

1 Findings	3
2 Introduction	3
2.1 Dataset	3
2.2 Data Cleaning	4
2.3 Validation	8
2.3.1 Random 1-Hold Out	8
2.3.2 Stratified 1-Hold Out	9
2.3.3 Random k-Fold	9
3 Random Forest	9
3.1 Split Metrics	9
3.1.1 Gini	9
3.1.2 Entropy	9
3.2 Experiments	10
4 Support Vector Machine	14
4.1 Experiments	14
5 Boosting	16
5.1 Adaboost	17
5.2 Random Forest	18
5.3 Experiment	18
5.4 Learning Rate and N Estimators	20
6 Appendix	21
6.1 Project Link	21
6.2 Code	21

List of Figures

1	Data visualization	3
2	Identify Noise Method	4
3	Identify Noise Method	5
4	Heat map	6
5	Heat Map and Filter Code	6
6	Feature 7	7
7	Pair plot	8
8	Gini Calculation	9
9	Entropy Calculation	9
10	Random Forest Code	10
11	Random Forest Raw Result	10
12	Depth Code	11
13	Depth Outcome	11
14	Split Code	12
15	Split Outcome	12
16	Class weight Code	12
17	Class Weight Outcome	13
18	Four error	13
19	SVM code	15
20	SVM Base Model Result	15
21	Boosting	17
22	Bagging	17
23	Adaboost Formula	18
24	Adaboost and Gradient Boost	18
25	Boosting Result	19
26	Bagging Result	20

1 Findings

In this experiment, first, I tried with different validation methods using base models. Then, I built my both model according to raw data. I continued with cleaning my Breast Cancer data. I made experience with both clean and raw data using both baby sited models. Then, I chose SVM as winner because it was giving more consistent results. Lastly, I applied both boosting and bagging methods. Boosting, slightly boosted the results, but bagging made it worse.

2 Introduction

In this take home exam assignment, our aim is to use methodologies learned earlier such as prepossessing and validation techniques with two classification methods, which are Support Vector Machine (SVM) and Random Forest. We are also aiming to apply bagging or boosting algorithm to the best.

2.1 Dataset

We are using Breast Cancer dataset for the take home exam assignment. It contains 2 classes, 569 instances with 30 features. We can understand by looking its shape. This is the shape of our data (569, 30).

In order to understand and get better perspective from our data. We project our data graph, and see how it looks.

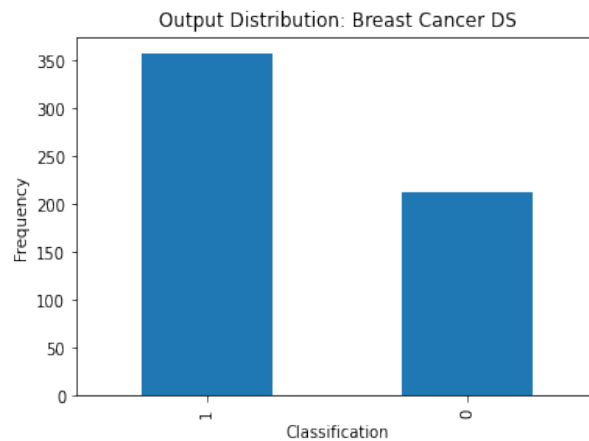


Figure 1: Data visualization

2.2 Data Cleaning

I conduct my experiments in two parts, one with raw data and one with cleaned data. In order to clean my data, I came up with some methods. First, I tried to identify the noisy data inside the data set.

```
def identify_noise(df):  
    noise = df[df.isnull().any(axis=1)].count()  
    total_noise = noise.sum()  
    print("{0} null values were found.".format(str(total_noise)))  
    if(total_noise > 0):  
        print(noise)  
    print("\n\nShowing all data types:\n\n")  
    print(df.dtypes)
```

Figure 2: Identify Noise Method

As a result, as we can see result below, there is no noisy data, and we can also see the data types of features.

```
0 null values were found.
```

```
Showing all data types:
```

```
0    float64
1    float64
2    float64
3    float64
4    float64
5    float64
6    float64
7    float64
8    float64
9    float64
10   float64
11   float64
12   float64
13   float64
14   float64
15   float64
16   float64
17   float64
18   float64
19   float64
20   float64
21   float64
22   float64
23   float64
24   float64
25   float64
26   float64
27   float64
28   float64
29   float64
dtype: object
```

Figure 3: Identify Noise Method

To see correlation in the data, I created a method, and visualize the data as a heat map. As we can see in below, the feature pairs that have a white color are highly correlated. Therefore, correlation will cause a problem for our models. For example, feature pairs (0, 2), (0, 3), and (2, 3) are all highly correlated. If we do not removed them, models will be highly correlated.

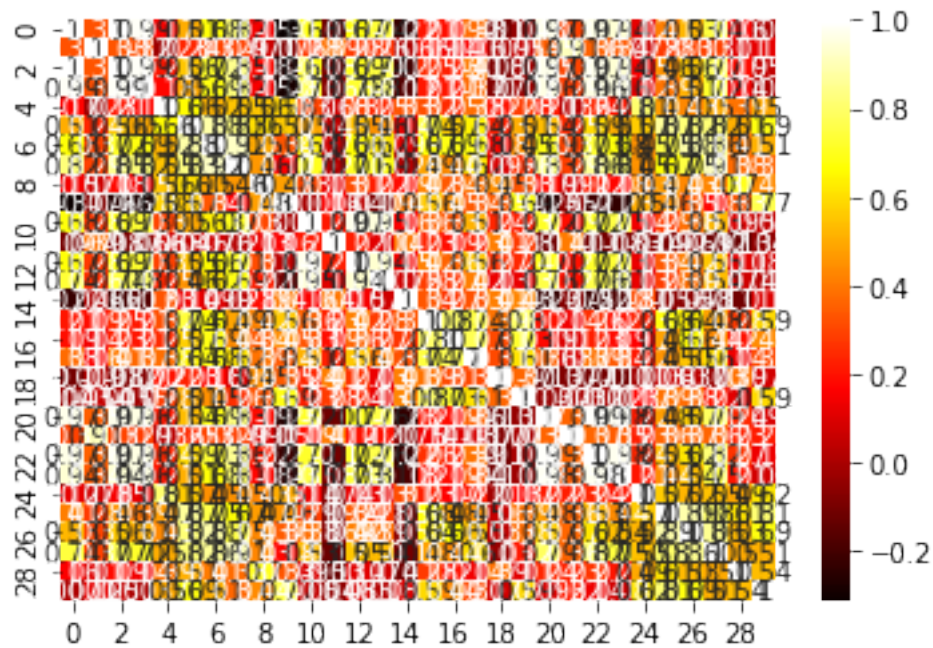


Figure 4: Heat map

```
def heat_map(df):
    fig, ax = plt.subplots()
    corr = df.corr()
    sns.heatmap(corr, annot=True, cmap='hot')
    plt.show()

def filter_features(data, bad_indices):
    # eliminate above column indices from the data and return new set
    filtered_data = np.delete(data, bad_indices, axis=1)

    return filtered_data
```

Figure 5: Heat Map and Filter Code

After done some cleaning using code above, again, we need to look graphs if there is a distinct difference between labels. For example, feature 7, looks like have distinct difference, so in the future, it will not be useful. Therefore, we need to remove it as well.

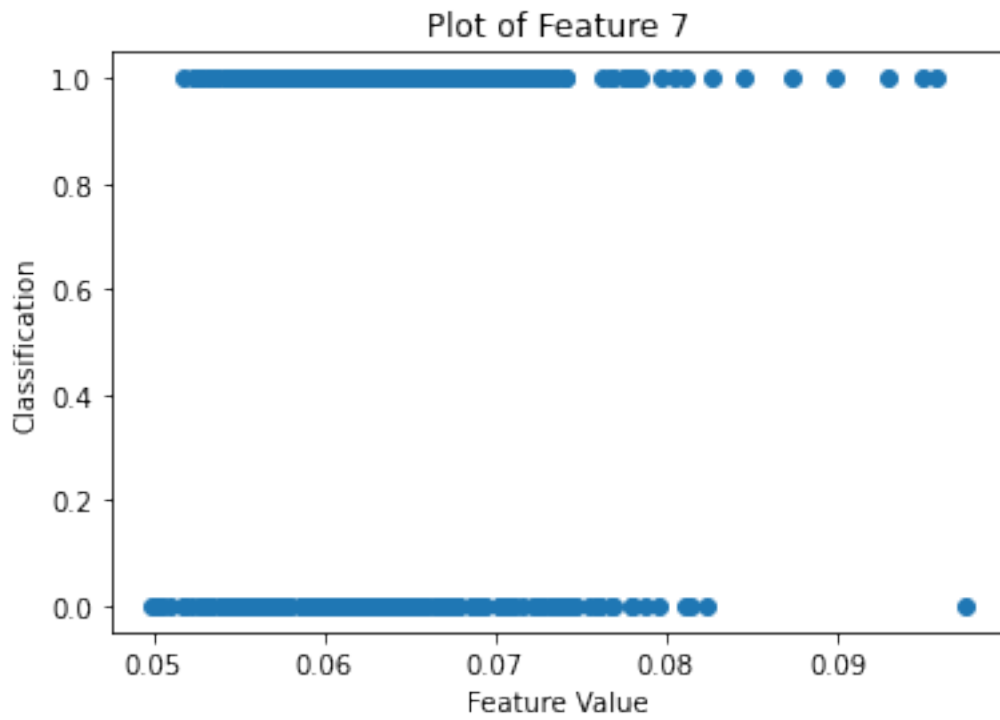


Figure 6: Feature 7

After the last cleaning, we have remained 14 features. Using the pairplot functionality, I display the data remain.

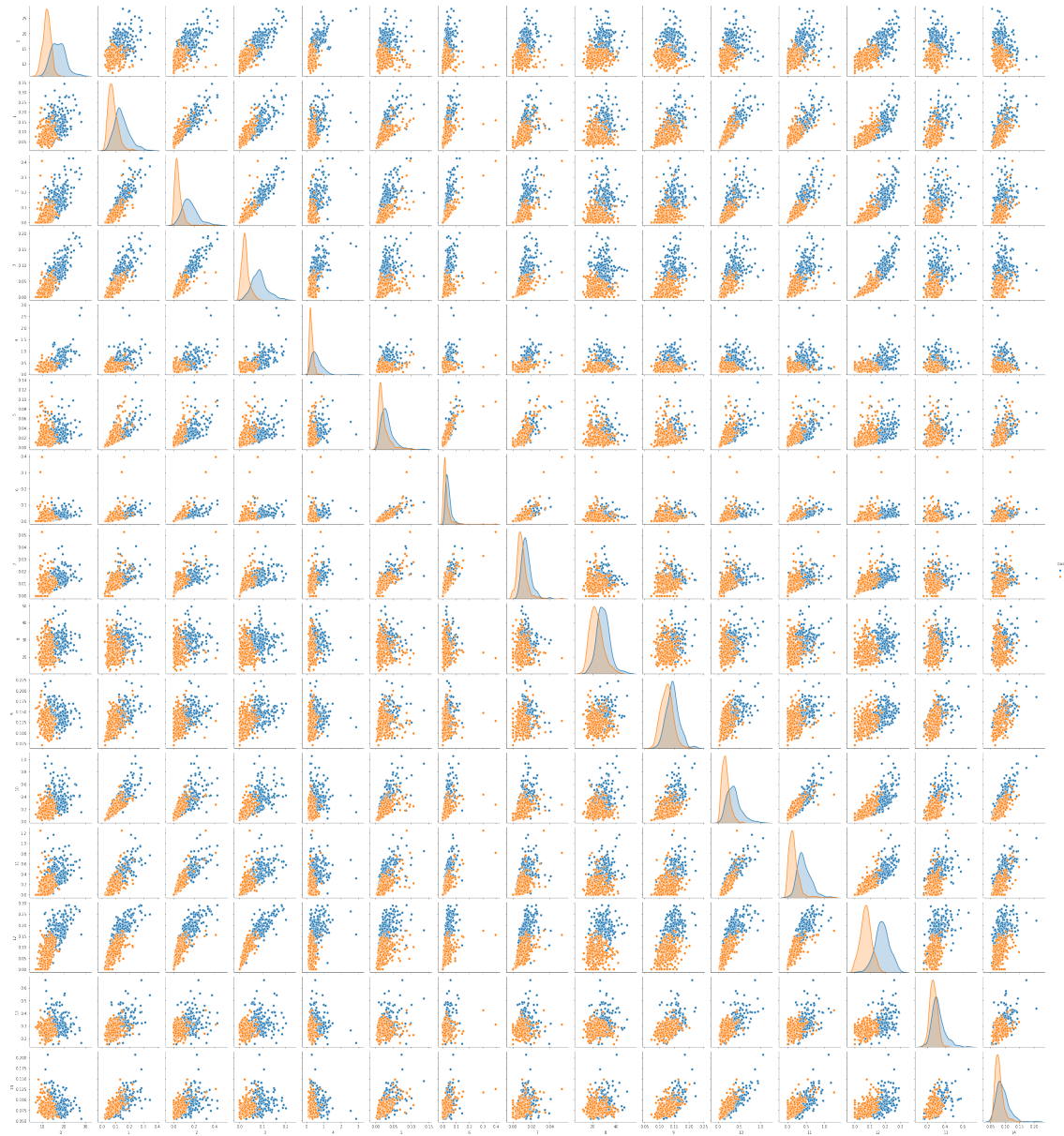


Figure 7: Pair plot

2.3 Validation

2.3.1 Random 1-Hold Out

Random 1-hold out is basically splitting up the dataset into a ‘train’ and ‘test’ set randomly. The training set is what the model is trained on, and the test is used to see performance of the model on unseen data.

2.3.2 Stratified 1-Hold Out

Stratified 1-hold out is splitting up the dataset into a ‘train’ and ‘test’ set so that each split has same percentage of samples of each targets as the complete set.

2.3.3 Random k-Fold

Random k-fold is a cross-validation technique which splits up the dataset into ‘k’ groups. One of the groups is used as the test set and the rest are used as the training set. The model is trained on the training set and scored on the test set. The process is repeated until each unique group as been used as the test set.

3 Random Forest

Random forests or random decision forests are an ensemble learning method for classification, regression and other tasks that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes or mean prediction of the individual trees. We will change it parameters such as weight, depth, and will try to find the best model.

3.1 Split Metrics

3.1.1 Gini

A Gini score gives an idea of how good a split is by how mixed the classes are in the two groups created by the split. A perfect separation results in a Gini score of 0, whereas the worst case split that results in 50/50 classes.

$$Gini = 1 - \sum_j p_j^2$$

Figure 8: Gini Calculation

3.1.2 Entropy

Entropy can be roughly thought of as how much variance the data has. It gives the homogeneity of a sample. If the sample is completely homogeneous the entropy is zero and if the sample is equally divided then it has entropy of one.

$$Entropy = - \sum_j p_j \log_2 p_j$$

Figure 9: Entropy Calculation

3.2 Experiments

First, we divided our data 70 percent as train data and 30 percent as test data. Then, we used different validation techniques on our train and test data to see which one is given the best result among them. For this purpose, we created a method called random_forest.

```
def random_forest(X_train, Y_train, forest_size, max_depth, criterion, min_samples_split, class_weight):
    rf = RandomForestClassifier(n_estimators=forest_size, oob_score=True, n_jobs=-1, max_depth=max_depth,
                               criterion=criterion, min_samples_split=min_samples_split, class_weight=class_weight, random_state=0)

    #5-Fold
    cv_result_rf_5 = cross_val_score(rf, X_train, Y_train, cv=5, scoring='accuracy')

    #10-Fold
    cv_result_rf_10 = cross_val_score(rf, X_train, Y_train, cv=10, scoring='accuracy')

    #Random One Holdout
    x_train, x_test, y_train, y_test_random = randomOneHoldout(X_train, Y_train)
    rf.fit(x_train, y_train)
    y_pred_rf_random = rf.predict(x_test)

    #Stratified One Holdout
    x_train, x_test, y_train, y_test_stratified = stratifiedOneHoldout(X_train, Y_train)
    rf.fit(x_train, y_train)
    y_pred_rf_stratified = rf.predict(x_test)

    print("Stratified One Hold Out Fold")
    print("Random Forest Accuracy: ", 1 - metrics.mean_squared_error(y_test_stratified, y_pred_rf_stratified))

    print("5 Fold")
    print("Random Forest Accuracy: ", cv_result_rf_5.mean())

    print("10 Fold")
    print("Random Forest Accuracy: ", cv_result_rf_10.mean())

    print("Random One Hold Out")
    print("Random Forest Accuracy: ", 1 - metrics.mean_squared_error(y_test_random, y_pred_rf_random))
```

Figure 10: Random Forest Code

Without manipulating the data, with base model of random forest, we got this accuracy. After the result, I decided to with ***Stratified One Hold Out***

```
Stratified One Hold Out Fold
Random Forest Accuracy: 0.975
5 Fold
Random Forest Accuracy: 0.959873417721519
10 Fold
Random Forest Accuracy: 0.9523717948717948
Random One Hold Out
Random Forest Accuracy: 0.975
```

Figure 11: Random Forest Raw Result

Then, I started to change its parameter individually. First, I started play with **depth** parameter. To do that, I created a method called **tuningDepth**.

```
def tuningDepth(x_train, x_test, y_train, y_test_stratified):
    max_depth_range = list(range(1, 10))
    max_depth_range.append(str("None"))

    for depth in max_depth_range:
        if(depth == "None"):
            clf = RandomForestClassifier(random_state = 0)
            clf.fit(x_train, y_train)
        else:
            clf = RandomForestClassifier(max_depth = depth, random_state = 0)
            clf.fit(x_train, y_train)

    accuracy = clf.score(x_test, y_test_stratified)*100
    print("Depth: ", depth, " Accuracy: ", accuracy)
```

Figure 12: Depth Code

After comparing our result with base model, as we can see, when the depth become higher, we can catch the base model. For the last model, I can use **7, 8, 9**

```
Stratified One Hold Out Fold
Random Forest Accuracy: 0.975

Depth: 1 Accuracy: 92.5
Depth: 2 Accuracy: 96.25
Depth: 3 Accuracy: 96.25
Depth: 4 Accuracy: 96.25
Depth: 5 Accuracy: 96.25
Depth: 6 Accuracy: 96.25
Depth: 7 Accuracy: 97.5
Depth: 8 Accuracy: 97.5
Depth: 9 Accuracy: 97.5
Depth: None Accuracy: 97.5
```

Figure 13: Depth Outcome

I continue to my experience with split strategy using **min_samples_split** and **criterion**. To do that I used code below.

```
def tuningSplit(x_train, x_test, y_train, y_test_stratified):
    criterion = ["gini", "entropy"]

    for i in criterion:
        clf = RandomForestClassifier(criterion = i, max_depth = 7, random_state = 0)
        clf.fit(x_train, y_train)

        accuracy = clf.score(x_test, y_test_stratified)*100
        print("Criterion: ", i, " Accuracy: ", accuracy)

    for i in range(2, 10):
        clf = RandomForestClassifier(max_depth = 7, min_samples_split = i, random_state = 0)
        clf.fit(x_train, y_train)

        accuracy = clf.score(x_test, y_test_stratified)*100
        print("min_samples_split: ", i, " Accuracy: ", accuracy)
```

Figure 14: Split Code

After comparing our result with base model, as we can see, when the **min_samples_split** become higher, we decreased our accuracy, and **criterion** did not change anything. For the last model, I can use **min_samples_split** as *2 or 3*

```
Stratified One Hold Out Fold
Random Forest Accuracy: 0.975

Criterion: gini Accuracy: 97.5
Criterion: entropy Accuracy: 97.5
min_samples_split: 2 Accuracy: 97.5
min_samples_split: 3 Accuracy: 97.5
min_samples_split: 4 Accuracy: 96.25
min_samples_split: 5 Accuracy: 96.25
min_samples_split: 6 Accuracy: 96.25
min_samples_split: 7 Accuracy: 96.25
min_samples_split: 8 Accuracy: 96.25
min_samples_split: 9 Accuracy: 96.25
```

Figure 15: Split Outcome

Lastly, individually, we changed the class_weight parameter. To do that, I used the code below.

```
def tuningClassWeight(x_train, x_test, y_train, y_test_stratified):
    # No class weight
    clf = RandomForestClassifier(max_depth = 7, random_state = 0)
    clf.fit(x_train, y_train)

    accuracy = clf.score(x_test, y_test_stratified)*100
    print("Class weight: None Accuracy: ", accuracy)

    # Balanced class weight
    clf = RandomForestClassifier(max_depth = 7, random_state = 0, class_weight = 'balanced')
    clf.fit(x_train, y_train)

    accuracy = clf.score(x_test, y_test_stratified)*100
    print("Class weight: Balanced Accuracy: ", accuracy)
```

Figure 16: Class weight Code

After comparing our result with base model, as we can see, when the **class_weight** does not change anything.

```
Stratified One Hold Out Fold
Random Forest Accuracy: 0.975

Class weight: None Accuracy: 97.5
Class weight: Balanced Accuracy: 97.5
```

Figure 17: Class Weight Outcome

In order to decide the final model, we conduct one more experiment. We combined all approaches mention above. This experiment happened inside **display-Accuracy** method After finding the best model, lastly, I tried to find optimal **n_estimator number**. Since the results are too much, I did not create a table, but I chose these parameters which are `n_estimator = 100`, `criterion = entropy`, `max_depth = 7`, `min_samples_split = 3` and `class_weight = "balanced"` to use in final model.

To check the correctness of the model, I checked the four errors written below:

```
def fourError(X, Y, model):
    X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.3, random_state=0, stratify=Y)

    Train_x, TrainDev_x, Train_y, TrainDev_y = train_test_split(X_train, Y_train, test_size=0.2, random_state=0, stratify=Y_train)
    Dev_x, Test_x, Dev_y, Test_y = train_test_split(X_test, Y_test, test_size=0.5, random_state=0, stratify=Y_test)

    model.fit(Train_x, Train_y)

    y_true, trainDev_pred = TrainDev_y, model.predict(TrainDev_x)

    print("Train-Train Dev, e1:", metrics.mean_squared_error(TrainDev_y, trainDev_pred), "\n")
    print("Accuracy: ", 1 - metrics.mean_squared_error(TrainDev_y, trainDev_pred))
    print(' \nClassification report\n' )
    print(classification_report(y_true, trainDev_pred))

    y_true, dev_pred = Dev_y, model.predict(Dev_x)
    print("Train-Dev, e2", metrics.mean_squared_error(Dev_y, dev_pred), "\n")
    print("Accuracy: ", 1 - metrics.mean_squared_error(Dev_y, dev_pred))
    print(' \nClassification report\n' )
    print(classification_report(y_true, dev_pred))

    y_true, test_pred = Test_y, model.predict(Test_x)
    print("Train-Test, e3: ", metrics.mean_squared_error(Test_y, test_pred), "\n")
    print("Accuracy: ", 1 - metrics.mean_squared_error(Test_y, test_pred))
    print(' \nClassification report\n' )
    print(classification_report(y_true, test_pred))

    y_true, devTest_pred = Y_test, model.predict(X_test)
    print("Train-(Dev+Test), e4: ", metrics.mean_squared_error(Y_test, devTest_pred), "\n")
    print("Accuracy: ", 1 - metrics.mean_squared_error(Y_test, devTest_pred))
    print(' \nClassification report\n' )
    print(classification_report(y_true, devTest_pred))
```

Figure 18: Four error

I tested the model with raw and cleaned data, and I saw that the results are okay, but the results are very inconsistent.

Train-Train Dev e1	Train-Dev e2	Train-Test e3	Train-(Dev+Test) e4
0.95	0.96470	0.91860	0.94152

Table 1: Four errors for raw data

Train-Train Dev e1	Train-Dev e2	Train-Test e3	Train-(Dev+Test) e4
0.925	0.98823	0.89534	0.94152

Table 2: Four errors for cleaned data

4 Support Vector Machine

The objective of the support vector machine algorithm is to find a hyper plane in an N-dimensional space(N — the number of features) that distinctly classifies the data points. To separate the two classes of data points, there are many possible hyper planes that could be chosen.

First, we divided our data 70 percent as train data and 30 percent as test data. Then, we used different validation techniques on our train and test data to see which one is given the best result among them. For this purpose, we created a method called SVM.

4.1 Experiments

To do our experiment, we created a method called SVM. we will apply our validation methods as well.

```
def svm(X_train, Y_train, kernel, weight, gamma):

    svm = SVC(C=1, kernel=kernel, degree=3, gamma=gamma, coef0=0.0, shrinking=True,
              probability=False, tol=0.001, cache_size=200, class_weight=weight,
              max_iter=-1, decision_function_shape="ovr", random_state = 0)

    #5-Fold
    cv_result_svm_5 = cross_val_score(svm, X_train, Y_train, cv=5, scoring='accuracy')

    #10-Fold
    cv_result_svm_10 = cross_val_score(svm, X_train, Y_train, cv=10, scoring='accuracy')

    #Random One Holdout
    x_train, x_test, y_train, y_test_random = randomOneHoldout(X_train, Y_train)
    svm.fit(x_train, y_train)
    y_pred_svm_random = svm.predict(x_test)

    #Stratified One Holdout
    x_train, x_test, y_train, y_test_stratified = stratifiedOneHoldout(X_train, Y_train)
    svm.fit(x_train, y_train)
    y_pred_svm_stratified = svm.predict(x_test)

    print("5 Fold")
    print("SVM Accuracy: ", cv_result_svm_5.mean())

    print("10 Fold")
    print("SVM Accuracy: ", cv_result_svm_10.mean())

    print("Random One Hold Out")
    print("SVM Accuracy: ", 1 - metrics.mean_squared_error(y_test_random, y_pred_svm_random))

    print("Stratified One Hold Out Fold")
    print("SVM Accuracy: ", 1 - metrics.mean_squared_error(y_test_stratified, y_pred_svm_stratified))
```

Figure 19: SVM code

We generally look our SVM model's result below, and we decided to move on with Stratified One Hold Out because it give higher accuracy. After we decided to go with Stratified One Hold Out , we started to feed our method with Kernel parameters. We used four different kernel parameters which are Linear, Poly, Rbf, Sigmoid.

```
5 Fold
SVM Accuracy: 0.8994936708860759
10 Fold
SVM Accuracy: 0.902051282051282
Random One Hold Out
SVM Accuracy: 0.9625
Stratified One Hold Out Fold
SVM Accuracy: 0.9625
```

File display

Figure 20: SVM Base Model Result

As we can see, we got the worst result by Sigmoid, and we got the best result by Linear. We continued with Linear. Then, we continued with class weight. In class_weight parameter, we have two option which are None and Balanced.

Because of the problem with colabs, I could not paste the results, but our result, 0.9875, stay same when we apply None parameter, but decreased when we applied Balanced. Finally, we made changes on gamma value. During our experiment, we figured out that Linear kernel only works with when gamma value "auto", and the result stayed same. As a result, we got quite good results, but I order to validate it, we need to look four error result. Also, I tried the model with cleaned data as well.

Train-Train Dev e1	Train-Dev e2	Train-Test e3	Train-(Dev+Test) e4
0.97802	0.94737	0.94737	0.94737

Table 3: Four errors for raw data

Train-Train Dev e1	Train-Dev e2	Train-Test e3	Train-(Dev+Test) e4
0.97802	0.94737	0.94737	0.94737

Table 4: Four errors for cleaned data

We got really good results, and they are most stable ones by far. We can apply boosting and bagging to final models.

5 Boosting

Boosting is one of most famous approaches and it produces an ensemble model that is in general less biased than the weak learners that compose it. Boosting methods work in the same spirit as bagging methods. We build a family of models that are aggregated to obtain a strong learner that performs better. However, unlike bagging that mainly aims at reducing variance, boosting is a technique that consists in fitting sequentially multiple weak learners in a very adaptative way.

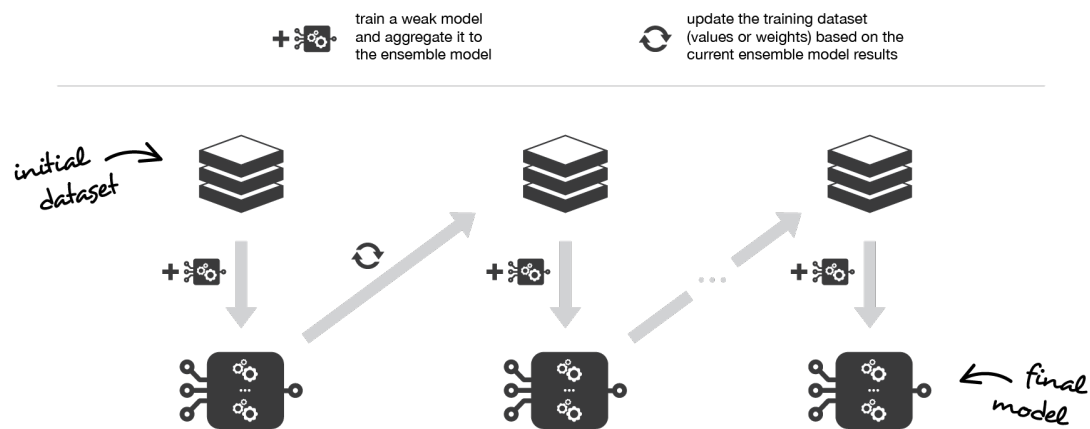


Figure 21: Boosting

Bagging, that often considers homogeneous weak learners, learns them independently from each other in parallel and combines them following some kind of deterministic averaging process.

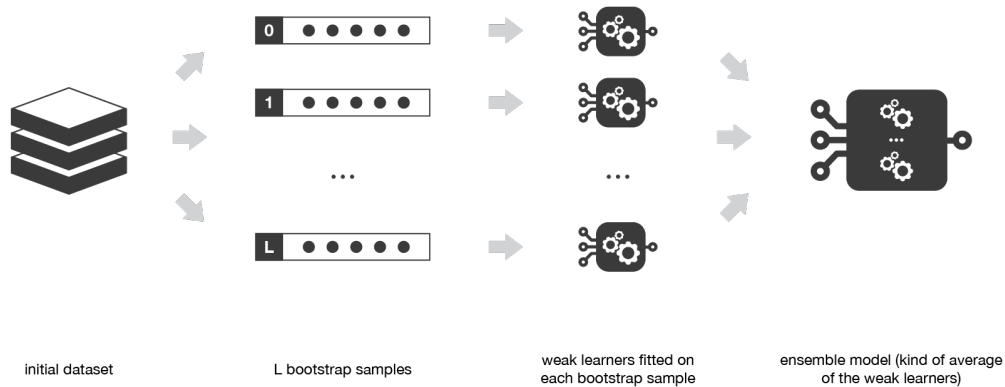


Figure 22: Bagging

5.1 Adaboost

In adaboosting (often called Adaptive Boost), we try to define our ensemble model as a weighted sum of L weak learners. In sklearn, we can feed our models to built-in adaboost method and we can get result.

$$s_L(.) = \sum_{l=1}^L c_l \times w_l(.) \quad \text{where } c_l\text{'s are coefficients and } w_l\text{'s are weak learners}$$

Figure 23: Adaboost Formula

5.2 Random Forest

In bagging, the random forest is one of the techniques for ensemble learning. The random forest approach where deep trees, fitted on bootstrap samples, are combined to produce an output with lower variance. However, random forests also use another trick to make the multiple fitted trees a bit less correlated with each others.

5.3 Experiment

For the experiment, we chose **SVM** comparing our results. We will apply boost and bagging methods to the both model. For boosting, we will compare it and change their two parameters which are learning_rate and n_estimators.

```
def AdaBoost(model, n_estimators, learning_rate, X_train, Y_train, X_test, Y_test):
    clf = AdaBoostClassifier(base_estimator = model, n_estimators= n_estimators, learning_rate=learning_rate, random_state=0)
    clf.fit(X_train, Y_train)
    clf.predict(X_test)
    return clf.score(X_train, Y_train)

def GradientBoost(n_estimators, learning_rate, X_train, Y_train, X_test, Y_test):
    clf = GradientBoostingClassifier(n_estimators= n_estimators, learning_rate=learning_rate, random_state=0)
    clf.fit(X_train, Y_train)
    clf.predict(X_test)
    return clf.score(X_train, Y_train)
```

Figure 24: Adaboost and Gradient Boost

For boosting, when the learners increased, we got lower accuracy. Also, when we changed the estimators, it did not affect the result. Overall, boosting increased our accuracy.

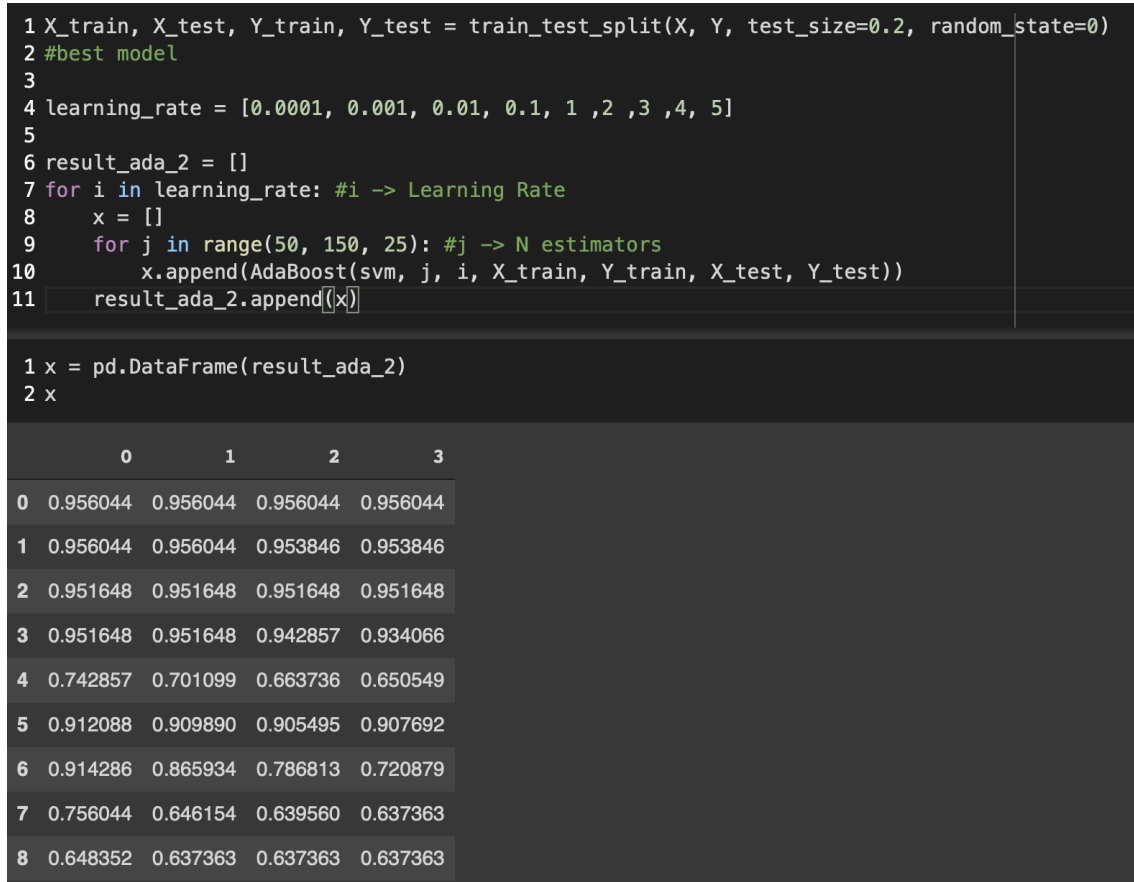


Figure 25: Boosting Result

For bagging, I only changed estimator number, and bagging gave pretty bad results.

```
1 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.3, random_state=0)
2 #best model
3
4
5 result_bagging = []
6 for j in range(50, 150, 25): #j -> N estimators
7     result_bagging.append(bagging(svm, j, X_train, Y_train, X_test, Y_test))
8
1 y = pd.DataFrame(result_bagging)
2 y
```

	0
0	0.736181
1	0.695980
2	0.665829
3	0.650754

Figure 26: Bagging Result

5.4 Learning Rate and N Estimators

Learning rate means that determines how much weak learners contribute to the weight of each iteration. Decreasing the learning rate makes the coefficients smaller, which reduces the amplitude of the sample_weights at each step. This translates into:

- Smaller variations of the weighted data points
- Fewer differences between the weak classifier decision boundaries

N estimators means that number of weak learners to train iteratively. Increasing the number of weak classifiers, increases the number of iterations, and allows the sample weights to gain greater amplitude. This translates into:

- More weak classifiers to combine at the end
- More variations in the decision boundaries of these classifiers

6 Appendix

6.1 Project Link

It will available after exam. https://github.com/kaankoken/random_forest_vs_svm

6.2 Code

```
0  # -*- coding: utf-8 -*-
1  """svm-vs-random-forest.ipynb
2
3  Automatically generated by Colaboratory.
4
5  Original file is located at
6      https://colab.research.google.com/github/kaankoken/
7          random_forest_vs_svm/blob/master/
8          svm_vs_random_forest.ipynb
9  """
10
11 from sklearn import metrics, datasets, preprocessing
12 import matplotlib.pyplot as plt
13 import numpy as np
14 import pandas as pd
15 import seaborn as sns
16 from sklearn.model_selection import train_test_split, KFold,
17     StratifiedKFold, cross_val_score
18
19 from sklearn.ensemble import RandomForestClassifier,
20     BaggingClassifier
21 from sklearn.preprocessing import StandardScaler
22 from sklearn.decomposition import PCA
23 from sklearn.svm import LinearSVC, SVC
24
25 from sklearn.ensemble import AdaBoostClassifier
26 from sklearn.metrics import classification_report
27
28 def randomOneHoldout(X_train, Y_train):
29     x_train, x_test, y_train, y_test = train_test_split(X_train,
30     Y_train, test_size=0.2,
31     random_state=0)
32
33     return x_train, x_test, y_train, y_test
34
35 def stratifiedOneHoldout(X_train, Y_train):
36     x_train, x_test, y_train, y_test = train_test_split(X_train,
37     Y_train, test_size=0.2,
38     random_state=0)
39
40     return x_train, x_test, y_train, y_test
```

```
32 def visualize_class(df):
    # look at the last column on data frame (the classification value
    # column)
34 df.iloc[:, -1].value_counts().plot(kind='bar')
    plt.title("Output Distribution: Breast Cancer DS")
36 plt.xlabel("Classification")
    plt.ylabel("Frequency")
38 plt.show()

40 def identify_noise(df):

42     noise = df[df.isnull().any(axis=1)].count()
    total_noise = noise.sum()
44     print("{0} null values were found.".format(str(total_noise)))
    if(total_noise > 0):
46         print(noise)
    print("\n\nShowing all data types:\n\n")
48     print(df.dtypes)

50 def heat_map(df):
    fig, ax = plt.subplots()
52     corr = df.corr()
    sns.heatmap(corr, annot=True, cmap='hot')
54     plt.show()

56 def filter_features(data, bad_indices):
    # eliminate above column indices from the data and return new set
58     filtered_data = np.delete(data, bad_indices, axis=1)

60     return filtered_data

62 def vis_all_feat(data, class_):
    for col_ind in range(data.shape[1]):
64         print("Viewing Feature #{0}".format(str(col_ind)))
        vis_single_feat(data, class_, col_ind)
66

def vis_single_feat(data, class_, ind):
68     # create graph of classification and feature values
    plt.figure(100) # display two plots on separate figures
70     df = pd.DataFrame(data)
    feat_vals = df.iloc[:, ind]
72     plt.scatter(feat_vals, class_)
    plt.title("Plot of Feature {0}".format(str(ind)))
74     plt.xlabel("Feature Value")
    plt.ylabel("Classification")
76

    # create bar graph of mean feature values for each classification
78     plt.figure(200)
    plt.title("Mean Values of Feature {0}".format(str(ind)))
80     plt.xlabel("Classification")
```

```

plt.ylabel("Mean Feature Value")
82 mean_df = pd.concat([df.iloc[:, ind], pd.Series(class_)], axis=1)
mean_df.columns = ["values", "classif"]
84 mean_df.groupby("classif", as_index=False)["values"].mean().loc[:,
                                                                "values"].plot(kind='bar')

86 plt.show()

88 def plot_pairplot(data, class_):
    data_df = pd.DataFrame(data)
90     # add classification so the plot can be colored by it
    data_df.loc[:, "classif"] = pd.Series(class_)
92     sns.pairplot(data_df, hue='classif')
    plt.show()

94 def fourError(X, Y, model):
96     X_train, X_test, Y_train, Y_test = train_test_split(X, Y,
                                                            test_size=0.2, random_state=0,
                                                            stratify=Y)

    Train_x, TrainDev_x, Train_y, TrainDev_y = train_test_split(
98         X_train, Y_train, test_size=0.2,
        random_state=0, stratify=Y_train)
    Dev_x, Test_x, Dev_y, Test_y = train_test_split(X_test, Y_test,
100         test_size=0.5, random_state=0,
        stratify=Y_test)

    model.fit(Train_x, Train_y)

102     y_true, trainDev_pred = TrainDev_y, model.predict(TrainDev_x)
104     print("Train-Train Dev,   e1:", metrics.mean_squared_error(
        TrainDev_y, trainDev_pred), "\n")
106     print("Accuracy: ", 1 - metrics.mean_squared_error(TrainDev_y,
        trainDev_pred))
    print( '\nClassification report\n' )
108     print(classification_report(y_true, trainDev_pred))

    y_true, dev_pred = Dev_y, model.predict(Dev_x)
110     print("Train-Dev,   e2", metrics.mean_squared_error(Dev_y,
        dev_pred), "\n")
112     print("Accuracy: ", 1 - metrics.mean_squared_error(Dev_y,
        dev_pred))
    print( '\nClassification report\n' )
114     print(classification_report(y_true, dev_pred))

    y_true, test_pred = Test_y, model.predict(Test_x)
116     print("Train-Test,   e3: ", metrics.mean_squared_error(Test_y,
        test_pred), "\n")

```



```
118     print("Accuracy: ", 1 - metrics.mean_squared_error(Test_y,
119                                                         test_pred))
120     print( '\nClassification report\n' )
121     print(classification_report(y_true, test_pred))
122
123     y_true, devTest_pred = Y_test, model.predict(X_test)
124     print("Train-(Dev+Test), e4: ", metrics.mean_squared_error(
125         Y_test, devTest_pred), "\n")
126     print("Accuracy: ", 1 - metrics.mean_squared_error(Y_test,
127                                                         devTest_pred))
128     print( '\nClassification report\n' )
129     print(classification_report(y_true, devTest_pred))
130
131 def svm(X_train, Y_train, kernel, weight, gamma):
132
133     svm = SVC(C=1, kernel=kernel, degree=3, gamma=gamma, coef0=0.0,
134               shrinking=True,
135               probability=False, tol=0.001, cache_size=200,
136               class_weight=weight,
137               max_iter=-1, decision_function_shape="ovr", random_state
138               = 0)
139
140     #5-Fold
141     cv_result_svm_5 = cross_val_score(svm, X_train, Y_train, cv=5,
142                                       scoring='accuracy')
143
144     #10-Fold
145     cv_result_svm_10 = cross_val_score(svm, X_train, Y_train, cv=10
146                                       , scoring='accuracy')
147
148     #Random One Holdout
149     x_train, x_test, y_train, y_test_random = randomOneHoldout(
150         X_train, Y_train)
151     svm.fit(x_train, y_train)
152     y_pred_svm_random = svm.predict(x_test)
153
154     #Stratified One Holdout
155     x_train, x_test, y_train, y_test_stratified =
156         stratifiedOneHoldout(X_train,
157                             Y_train)
158     svm.fit(x_train, y_train)
159     y_pred_svm_stratified = svm.predict(x_test)
160
161     print("5 Fold")
162     print("SVM Accuracy: ", cv_result_svm_5.mean())
163
164     print("10 Fold")
165     print("SVM Accuracy: ", cv_result_svm_10.mean())
166
167     print("Random One Hold Out")
```

```
print("SVM Accuracy: ", 1 - metrics.mean_squared_error(
    y_test_random, y_pred_svm_random)
)

158 print("Stratified One Hold Out Fold")
160 print("SVM Accuracy: ", 1 - metrics.mean_squared_error(
    y_test_stratified,
    y_pred_svm_stratified))

162 def random_forest(X_train, Y_train, forest_size, max_depth,
    criterion, min_samples_split,
    class_weight):
    rf = RandomForestClassifier(n_estimators=forest_size, oob_score
    =True, n_jobs=-1, max_depth=
    max_depth, criterion=criterion,
    min_samples_split =
    min_samples_split, class_weight=
    class_weight, random_state=0)

164
166 #5-Fold
    cv_result_rf_5 = cross_val_score(rf, X_train, Y_train, cv=5,
    scoring='accuracy')

168 #10-Fold
    cv_result_rf_10 = cross_val_score(rf, X_train, Y_train, cv=10,
    scoring='accuracy')

170
172 #Random One Holdout
    x_train, x_test, y_train, y_test_random = randomOneHoldout(
    X_train, Y_train)
    rf.fit(x_train, y_train)
    y_pred_rf_random = rf.predict(x_test)

174
176 #Stratified One Holdout
    x_train, x_test, y_train, y_test_stratified =
    stratifiedOneHoldout(X_train,
    Y_train)
    rf.fit(x_train, y_train)
    y_pred_rf_stratified = rf.predict(x_test)

178
180 print("Stratified One Hold Out Fold")
182 print("Random Forest Accuracy: ", 1 - metrics.
    mean_squared_error(
    y_test_stratified,
    y_pred_rf_stratified))

184 print("5 Fold")
    print("Random Forest Accuracy: ", cv_result_rf_5.mean())

186 print("10 Fold")
```

```
188     print("Random Forest Accuracy: ", cv_result_rf_10.mean())
190
191     print("Random One Hold Out")
192     print("Random Forest Accuracy: ", 1 - metrics.
193           mean_squared_error(y_test_random,
194                             y_pred_rf_random))
195
196 def tuningDepth(x_train, x_test, y_train, y_test_stratified):
197
198     max_depth_range = list(range(1, 10))
199     max_depth_range.append(str("None"))
200
201     for depth in max_depth_range:
202         if(depth == "None"):
203             clf = RandomForestClassifier(random_state = 0)
204             clf.fit(x_train, y_train)
205         else:
206             clf = RandomForestClassifier(max_depth = depth,
207                                       random_state = 0)
208             clf.fit(x_train, y_train)
209
210         accuracy = clf.score(x_test, y_test_stratified)*100
211         print("Depth: ", depth, " Accuracy: ", accuracy)
212
213 def tuningSplit(x_train, x_test, y_train, y_test_stratified):
214     criterion = ["gini", "entropy"]
215
216     for i in criterion:
217         clf = RandomForestClassifier(criterion = i, max_depth = 7,
218                                   random_state = 0)
219         clf.fit(x_train, y_train)
220
221         accuracy = clf.score(x_test, y_test_stratified)*100
222         print("Criterion: ", i, " Accuracy: ", accuracy)
223
224     for i in range(2, 10):
225         clf = RandomForestClassifier(max_depth = 7,
226                                   min_samples_split = i,
227                                   random_state = 0)
228         clf.fit(x_train, y_train)
229
230         accuracy = clf.score(x_test, y_test_stratified)*100
231         print("min_samples_split: ", i, " Accuracy: ", accuracy)
232
233 def tuningClassWeight(x_train, x_test, y_train, y_test_stratified):
234     # No class weight
235     clf = RandomForestClassifier(max_depth = 7, random_state = 0)
236     clf.fit(x_train, y_train)
237
238     accuracy = clf.score(x_test, y_test_stratified)*100
```

```

232     print("Class weight: None           Accuracy: ", accuracy)
234
234     # Balanced class weight
234     clf = RandomForestClassifier(max_depth = 7, random_state = 0,
234                                class_weight = 'balanced')
236     clf.fit(x_train, y_train)
238
238     accuracy = clf.score(x_test, y_test_stratified)*100
238     print("Class weight: Balanced       Accuracy: ", accuracy)
240
240     def AdaBoost(model, n_estimators, learning_rate, X_train, Y_train,
240                  X_test, Y_test):
242         clf = BaggingClassifier(base_estimator = model, n_estimators=
242                                n_estimators, learning_rate=
242                                learning_rate, random_state=0)
244
244         clf.fit(X_train, Y_train)
244         clf.predict(X_test)
244         return clf.score(X_train, Y_train)
246
246     def bagging(model, n_estimators, X_train, Y_train, X_test, Y_test):
248         clf = AdaBoostClassifier(base_estimator = model, n_estimators=
248                                n_estimators, random_state=0)
250
250         clf.fit(X_train, Y_train)
250         clf.predict(X_test)
250         return clf.score(X_train, Y_train)
252
252     def displayAccuracy(X, Y):
254         X_train, X_test, Y_train, Y_test = train_test_split(X, Y,
254                                                             test_size=0.3, random_state=0)
256
256         kernel = ["linear", "rbf", "poly", "sigmoid"]
256         weight = [None, "balanced"]
258         gamma = ["auto", "scale"]
260
260         for i in kernel:
260             for j in weight:
262                 for k in gamma:
262                     if i != "linear":
264                         print("Kernel: {} - Weight: {} - Gamma: {}".
264                               format(str(i), j, k))
264                         svm(X_train, Y_train, i, j, k)
266                     else:
266                         print("Kernel: {} - Weight: {} - Gamma: {}".
266                               format(str(i), j, "auto"))
268                         svm(X_train, Y_train, i, j, k)
268
268                 print()
270
270         combined approach
272         for depth in max_depth_range:
272             for c in criterion:

```

```

274         for i in range(2, 10):
275             for w in weight:
276                 print("Depth: {}, Criterion: {}, Min Split {},
                        Weight: {}".format(depth, c, i, w
                        ))
                        random_forest(X_train, Y_train, 100, depth, c,
277                                     i, w)
278         print("\n")
279
280     forest_size = [25, 50, 75, 100, 125, 150, 175, 200]
281
282     for s in forest_size:
283         random_forest(X_train, Y_train, s, 7, "entropy", 3, "
                        balanced")
284
285     def compareBaseModels(X, Y):
286         X_train, X_test, Y_train, Y_test = train_test_split(X, Y,
                        test_size=0.3, random_state=0)
287
288         svm(X_train, Y_train, "rbf", None, "scale")
289         print("")
290         random_forest(X_train, Y_train, 100, None, "gini", 2, None)
291
292         x_train, x_test, y_train, y_test_stratified =
                        stratifiedOneHoldout(X_train,
                        Y_train)
293
294         #Individual tuning
295         print("\n")
296         tuningDepth(x_train, x_test, y_train, y_test_stratified);
297         tuningSplit(x_train, x_test, y_train, y_test_stratified)
298         tuningClassWeight( x_train, x_test, y_train, y_test_stratified)
299
300     if __name__ == '__main__':
301
302         breast_cancer = datasets.load_breast_cancer()
303         X = breast_cancer.data
304         Y = breast_cancer.target
305
306         #Shape of the data
307         print(X.shape, end="\n")
308
309         feature = pd.DataFrame(Y)
310         df = pd.DataFrame(X)
311         rf = RandomForestClassifier(n_estimators=100, oob_score=True,
                        n_jobs=-1, max_depth=7, criterion
                        ="entropy", min_samples_split = 3
                        , class_weight="balanced",
                        random_state=0)
312         svm = SVC(C=1, kernel="linear", degree=3, gamma="scale", coef0=
                        0.0, shrinking=True,

```

```
312         probability=True, tol=0.001, cache_size=200, class_weight
           =None,
           max_iter=-1, decision_function_shape="ovr", random_state
           = 0)
314 fourError(X, Y, svm)
           visualize_class(feature)
316 compareBaseModels(X, Y)
           displayAccuracy(X, Y)
318
           identify_noise(df)
320 heat_map(df)
322 #filter strongly correlated features - can see which ones in
           correlation map
           X = filter_features(X, [2, 3, 20, 22, 23, 12, 13])
324
           vis_all_feat(X, Y)
326 X = filter_features(X, [1, 2, 6, 7, 9, 10, 14, 15])
           print("Cleaned data")
328 fourError(X, Y, svm)
           #remaining features
330 plot_pairplot(X, Y)
332 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size
           =0.3, random_state=0)
           #best model
334
           learning_rate = [0.0001, 0.001, 0.01, 0.1, 1 ,2 ,3 ,4, 5]
336
           result_ada_2 = []
338 for i in learning_rate: #i -> Learning Rate
           x = []
340     for j in range(50, 150, 25): #j -> N estimators
           x.append(AdaBoost(svm, j, i, X_train, Y_train, X_test,
                           Y_test))
342     result_ada_2.append(x)
344 x = pd.DataFrame(result_ada_2)
           x
346
           X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size
           =0.3, random_state=0)
348 #best model
350
           result_bagging = []
352 for j in range(50, 150, 25): #j -> N estimators
           result_bagging.append(bagging(svm, j, X_train, Y_train, X_test,
                           Y_test))
354
```

```
356 y = pd.DataFrame(result_bagging)
    y
```