

## Tox21 project report

Arash Sajjadi

Applied Machine Learning

Shahid Beheshti University

February 22, 2022



**Keywords:** Machine learning, Machine learning algorithms, K-NN, K-nearest neighbors, Support vector machine, svm, Gridsearchcv, Boosting, Bagging, AdaBoost, Ensemble, Random forest, Decision Tree, Python, MLP, Balanced accuracy, Confusion matrix

### ABSTRACT

The raw Tox21 data in this report are related to the toxicity of various chemicals. The number of raw data is 7830. Labels are 12-dimensional vectors, each dimension indicating whether the chemical molecule is toxic or not from the point of view of that dimension. My task is to find a function that can compute the prediction of this 12-dimensional vector. [1] book, I tried to figure out the relation between the data and whether they were effective in treating HIV or not. I hope you enjoy the following content.

### CONTENTS

<b>Contents</b>	<b>1</b>
<b>1</b> Introduction	<b>1</b>
<b>2</b> Description of raw data	<b>2</b>
2.1 Labels distribution . . . . .	2
2.2 Missing data . . . . .	2
2.3 Algebraic dimension of features . . . . .	2
<b>3</b> Data separation	<b>3</b>
<b>4</b> dimensionality reduction	<b>3</b>
<b>5</b> Correlation of labels	<b>3</b>
<b>6</b> Algorithms	<b>3</b>
6.1 K Nearest Neighbors . . . . .	4
6.2 Decision Tree . . . . .	6
6.3 Random forest . . . . .	6
6.4 XGBoost . . . . .	7
6.5 C-Support Vector Classification . . . . .	8
6.6 AdaBoost . . . . .	9
6.7 Bagging . . . . .	10
6.8 MLP . . . . .	11
<b>7</b> Evaluation	<b>12</b>
7.1 1st dimension of the labels vector . . . . .	12
7.2 2nd dimension of the labels vector . . . . .	12
7.3 3rd dimension of the labels vector . . . . .	12
7.4 4th dimension of the labels vector . . . . .	13
7.5 5th dimension of the labels vector . . . . .	14
7.6 6th dimension of the labels vector . . . . .	14
7.7 7th dimension of the labels vector . . . . .	14
7.8 8th dimension of the labels vector . . . . .	14
7.9 9th dimension of the labels vector . . . . .	15

7.10 10th dimension of the labels vector . . . . .	16
7.11 11th dimension of the labels vector . . . . .	17
7.12 12th dimension of the labels vector . . . . .	17
<b>8</b> conclusion	<b>18</b>
8.1 Initial conclusion . . . . .	18
8.2 Final project conclusion . . . . .	21
<b>References</b>	<b>23</b>

### PRE-INTRODUCTION

Toxicity is an adverse health effect caused by a drug. The adverse effect can range from a minor unpleasant side effect to a major threat to the life quality, health, or survival of the patient. [2]

You may be would like to know why it is essential for us to check the toxicity of chemicals. Knowing whether a chemical can cause cancer, allergic reactions, or abnormalities in unborn children is vital to human health. The process of discovering this information is known as the assessment of hazard.

### 1 INTRODUCTION

My data contains two files named `tox21_global_cdf_rdkit` and `tox21`. The first one is for the properties of each data, and the second one is for the labels. Also, the chemical formula of molecules is written in strings format, called Smiles format. For example, the eighth data refers to the chemical molecule in Figure 1.[3]

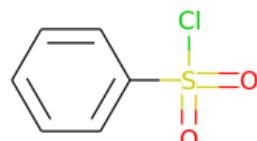


Figure 1: The chemical formula of the eighth data

The features file contains a small number of empty variables. For this reason, I can delete the rows containing the empty variable from the beginning. However, the labels file contains many vectors that have at least one empty variable. If I want to delete all vectors containing at least one empty variable, only 3054<sup>1</sup> data will remain. To train a good

model, I have to think of a better idea. I will mention this idea in the DESCRIPTION OF RAW DATA section.

## 2 DESCRIPTION OF RAW DATA

### 2.1 Labels distribution

As I mentioned before, each data label is a 12-dimensional vector. Each dimension informs us of the toxicity of that chemical in binary form. The labels of each dimension of data have an imbalanced distribution. Figure 2 shows the distribution of each dimension for us. Of course, these charts have been prepared after the clever removal of missing data, which I will explain in more detail in the next subsection.

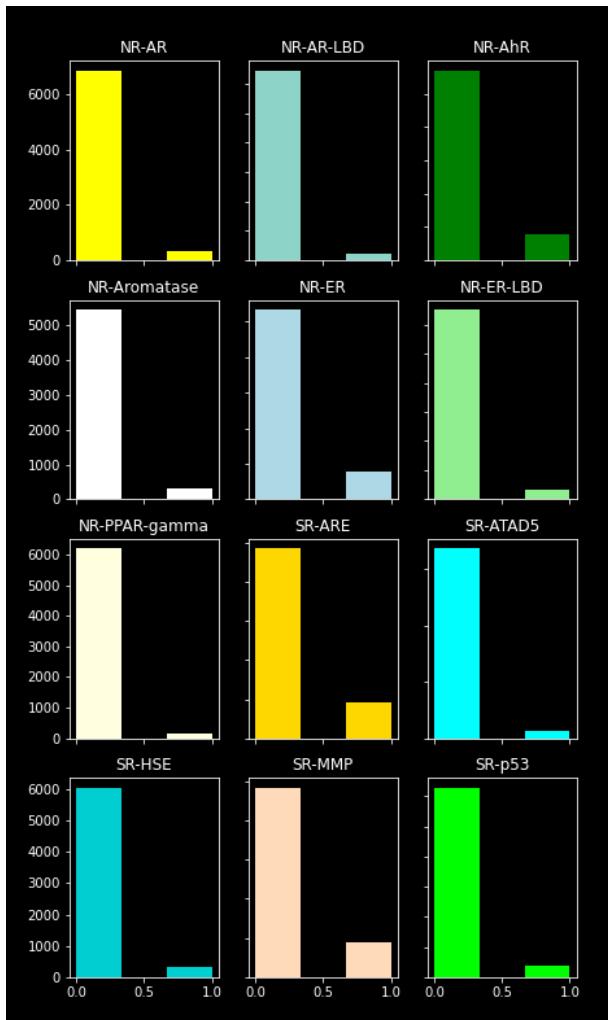


Figure 2: Labels distribution [4]

### 2.2 Missing data

I have probably been able to describe the challenge that missing data poses to us so far. To deal with this problem, instead of eliminating the rows containing empty data, we examine each dimension separately and

remove only the empty data related to that dimension.

$$\begin{bmatrix} \text{NR-AR} & \text{NR-AR-LBD} & \dots & \text{SR-p53} \\ \text{NaN} & \text{NaN} & \dots & \text{NaN} \\ \text{NaN} & \text{NaN} & \dots & \text{NaN} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & \dots & 0 \end{bmatrix} \mapsto \begin{bmatrix} \text{NR-AR} \\ \text{NaN} \\ 0 \\ \vdots \\ 1 \end{bmatrix}, \begin{bmatrix} \text{NR-AR-LBD} \\ \text{NaN} \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} \text{SR-p53} \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

As you can see in Figure 2, none of the labels are distributed in a balanced manner. The data imbalance causes us to have other problems, for which a solution has been predicted. For example, in the case of separating the data set into the training set and test set, we must use the stratify data option when separating the data.[5] you can see more detail in Listing 1

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import warnings
5 warnings.simplefilter('ignore')
6 #importing packages
7 initial_Features=pd.read_csv('tox21_global_cdf_rdkit.
8   csv')
9 initial_dataset=pd.read_csv('tox21.csv')
10 #import data
11 initial_Features=initial_Features.loc[:,,
12   initial_Features.apply(pd.Series.nunique) != 1]
13 initial_dataset=initial_dataset.iloc[initial_Features.
14   dropna().index]
15 #algebraic dimension reduction
16 initial_dataset=initial_dataset.reset_index()
17 initial_Features=initial_Features.dropna()
18 initial_Features=initial_Features.reset_index()
19 index_array=[]
20 for i in np.arange(1,13):
21   index_array.append(initial_dataset.iloc[:,i+1].
22     dropna().index)
23 def label_ith(i):
24   return pd.DataFrame(data=initial_dataset.iloc[
25     index_array[i]].iloc[:,i+2])
26 def Feature_ith(i):
27   return initial_Features.iloc[index_array[i]].drop(
28     'index',axis=1)
29 #define functions to sync Nan removed labels with
30   features
31 from sklearn.model_selection import train_test_split
32 X_training_data=[]
33 X_test=[]
34 y_training_data=[]
35 y_test=[]
36 for i in np.arange(0,12):
37   X_training_data_tmp, X_test_tmp,
38   y_training_data_tmp, y_test_tmp =train_test_split(
39     Feature_ith(i),label_ith(i), stratify=label_ith(i),
40     test_size=0.10,random_state=1234)
41   X_training_data.append(X_training_data_tmp)
42   X_test.append(X_test_tmp)
43   y_training_data.append(y_training_data_tmp)
44   y_test.append(y_test_tmp)
45 #split, sync data as we saw in the matrix before.

```

Listing 1: Import packages Deal with missing data Data separation [6]  
[7]

### 2.3 Algebraic dimension of features

Raw data is 200 dimensions. I quickly noticed that some of the columns of the matrix feature show a fixed number. By removing these columns, we have 190 algebraic dimensions.[Refer to Listing 1] Also, the data was normalized from the beginning. 201st column was related to the data normalization report, which I did not consider as a separate dimension from the beginning.

### 3 DATA SEPARATION

It is clear that in machine learning projects, one of the first tasks is to separate the test set. I must do this before any data analysis. If the data separation is to be done after PCA<sup>2</sup> or anything like that, we have practically used the test set data to find the best directions, which is a mistake. As a result, data separation generally must be done at the beginning of machine learning projects.<sup>3</sup> I also intended to allocate 10% of the data volume to the test set.

As Listing 1, there are functions called `label_iith` and `Feature_iith`, which give us labels and data synced with nonempty labels, respectively. Also, in each cell, the `X_training_data`, `y_training_data`, `X_test`, and `y_test` arrays, as their names, contain information synced to the same label dimension.

### 4 DIMENSIONALITY REDUCTION

High-dimensional data has some peculiar characteristics, some of which are counterintuitive. For example, in high dimensions the center of the space is devoid of points, with most of the points being scattered along the surface of the space or in the corners. There is also an apparent proliferation of orthogonal axes. Consequently, high-dimensional data can cause problems for data mining and analysis, although in some cases, high-dimensionality can help, for example, for nonlinear classification. Nevertheless, it is essential to check whether the dimensionality can be reduced while preserving the crucial properties of the full data matrix. One of the critical methods for dimension reduction is PCA.

Principal Component Analysis (PCA) is a technique that seeks an  $r$ -dimensional basis that best captures the variance in the data. The direction with the largest projected variance is called the first principal component. The orthogonal direction that captures the second largest projected variance is called the second principal component, and so on. As we shall see, the direction that maximizes the variance is also the one that minimizes the mean squared error [8].

Since I have prepared arrays such as `X_training_data`, etc., for each dimension of the label vector, applying the PCA algorithm separately for each member of these arrays is better. You can also see the results of applying this algorithm in Listing 2. It is also evident that the test set information for the PCA algorithm has not been used in any way. I found that my data with 71 principal components can maintain more than 95% of the data variance in each of the 12 data sets.

```

1 from sklearn.decomposition import PCA
2 X_training_data_pca=[]
3 X_test_pca=[]
4 for i in np.arange(0,12):
5     pca = PCA(n_components=71)
6     principalComponents = pca.fit_transform(
7         X_training_data[i])
8     X_training_data_PCA_tmp = pd.DataFrame(data =
9         principalComponents)
10    X_training_data_pca.append(X_training_data_PCA_tmp)
11
12    X_test_pca_tmp=pd.DataFrame(data=pca.transform(
13        X_test[i]))
14    X_test_pca.append(X_test_pca_tmp)
15    print('PCA with 71 principal components retains',
16        np.sum(pca.explained_variance_ratio_)*100,'% of
17        data VAR. (It is related for',i,'th label)')

```

Listing 2: PCA [5]

```

4 PCA with 71 principal components retains 95.24899205532745 % of data VAR. (It is
related for 3 th label)
5 PCA with 71 principal components retains 95.19568949095981 % of data VAR. (It is
related for 4 th label)
6 PCA with 71 principal components retains 95.18825061528679 % of data VAR. (It is
related for 5 th label)
7 PCA with 71 principal components retains 95.1389446393015 % of data VAR. (It is
related for 6 th label)
8 PCA with 71 principal components retains 95.21618305641543 % of data VAR. (It is
related for 7 th label)
9 PCA with 71 principal components retains 95.18018989506741 % of data VAR. (It is
related for 8 th label)
10 PCA with 71 principal components retains 95.16845227328498 % of data VAR. (It is
related for 9 th label)
11 PCA with 71 principal components retains 95.22548243540908 % of data VAR. (It is
related for 10 th label)
12 PCA with 71 principal components retains 95.17820123387547 % of data VAR. (It is
related for 11 th label)

```

Listing 3: PCA results

### 5 CORRELATION OF LABELS

Before running the algorithms, I was personally asked how correlated the different dimensions of the labels were. I want to talk here as someone who has seen the end of the project. We were much more successful in predicting some labels in proportion to some others. So if we could find a high correlation between the labels, it would be possible to use the result of more successful predictions for less accurate predictions. But as you can see in Figures 3 and 4, there is not much correlation between the labels.

Is it not harmful to explain the difference between Figures 3 and 4. Figure 3 shows the correlation between the raw data. We know that the raw data contained a significant amount of empty data. Figure 4 shows the correlation of non-empty data labels. What is clear is that the maximum correlation between the data is 0.551974651, which clearly indicates that different data labels are not going to help us predict other labels.

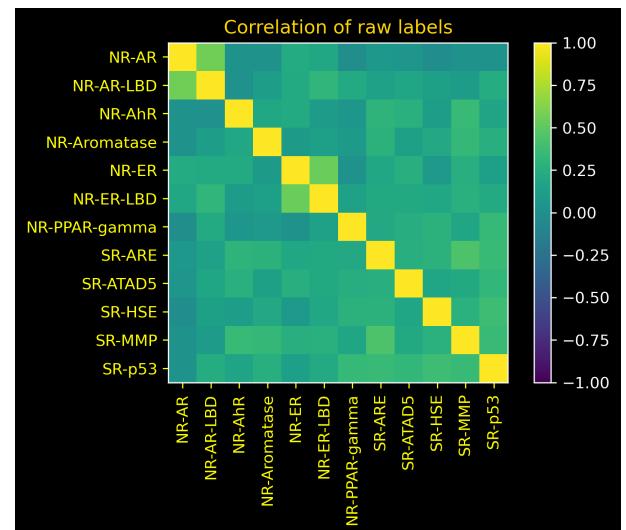


Figure 3: Correlation of the raw labels [4]

### 6 ALGORITHMS

In this section, I will execute each classification algorithm for each dimension of the label vectors with the help of cross-validation<sup>4</sup> and share the results with you. One of my efforts during this project has been to keep the code short for better readability. I hope you enjoy these beautiful algorithms.

<sup>1</sup> PCA with 71 principal components retains 95.14902477458284 % of data VAR. (It is related for 0 th label)  
<sup>2</sup> PCA with 71 principal components retains 95.18528509261213 % of data VAR. (It is related for 1 th label)  
<sup>3</sup> PCA with 71 principal components retains 95.23516716185081 % of data VAR. (It is related for 2 th label)

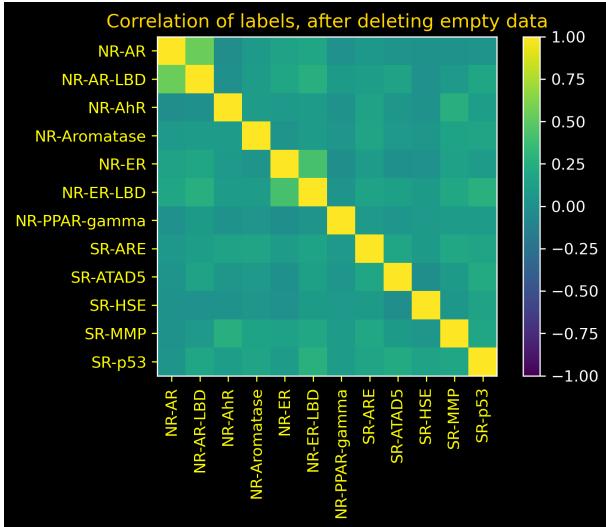


Figure 4: Correlation of labels, after deleting empty data [4]

### Scoring meter

Given the nature of the imbalanced data, I decided to choose the balanced accuracy criterion to evaluate the results of my algorithms. This accuracy measurement is defined as follows.

$$\text{Balanced Accuracy} = \frac{1}{2} \left( \frac{\text{TP}}{\text{P}} + \frac{\text{TN}}{\text{N}} \right)$$

Of course, if necessary, I will use the confusion matrix [Refer to Figure 5] to analyze the algorithms.

		True Class	
		Negative	Positive
Predicted Class	Negative	TN	FP
	Positive	FN	TP

Figure 5: confusion matrix[9]<sup>6</sup>

### 6.1 K Nearest Neighbors

Nearest Neighbor algorithms are among the simplest of all machine learning algorithms. The idea is to memorize the training set and then to predict the label of any new instance based on the labels of its closest neighbors in the training set. The rationale behind such a method is based on the assumption that the features that are used to describe the domain points are relevant to their labelings in a way that makes close-by points likely to have the same label. Furthermore, in some situations, even when the training set is immense, finding the nearest neighbor can be done extremely fast <sup>7</sup>. [1] You can see the execution of KNN on 77-dimensional data (after PCA application) in Listing 4.

```

1 from sklearn.model_selection import GridSearchCV
2 from sklearn.metrics import confusion_matrix,
3     plot_confusion_matrix
4 print("KNeighborsClassifier:")
5 for i in np.arange(0,12,1):
6     from sklearn.neighbors import KNeighborsClassifier
7     knn = KNeighborsClassifier()
8     param_grid={'n_neighbors':np.arange(1,11,2),
9                 'metric':['euclidean','cosine','manhattan']}
10    grid=GridSearchCV(knn,param_grid=param_grid, cv=10,
11                      scoring='balanced_accuracy', return_train_score=False)
12    grid.fit(X_training_data_pca[i].to_numpy(),
13              y_training_data[i])
14    print("----This information is for the ",i,"th
15          label.----")
16    print("-----Including PCA
17          -----")
18    print("best mean cv score= ",grid.best_score_)
19    print("best parameters= ",grid.best_params_)
20    y_pred_knn_pca=grid.best_estimator_.predict(
21        X_test_pca[i])
22    print(confusion_matrix(y_test[i],np.transpose(np.
23        matrix(y_pred_knn_pca))))
24    print(grid.score(X_test_pca[i],y_test[i]))

```

Listing 4: GridSearchCV-KNN-PCA [5]

```

1 KNeighborsClassifier:
2 ---This information is for the 0 th label.---
3 -----Including PCA.-----
4 best mean cv score= 0.7235540265998088
5 best parameters= {'metric': 'manhattan', 'n_neighbors': 7}
6 [[682 4]
7  [ 15 16]]
8 0.7551490642339885
9 ---This information is for the 1 th label.---
10 -----Including PCA.-----
11 best mean cv score= 0.7850023245692129
12 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
13 [[637 9]
14  [ 6 17]]
15 0.8625992731188585
16 ---This information is for the 2 th label.---
17 -----Including PCA.-----
18 best mean cv score= 0.7259460692370788
19 best parameters= {'metric': 'cosine', 'n_neighbors': 1}
20 [[544 27]
21  [ 33 44]]
22 0.7620715536652489
23 ---This information is for the 3 th label.---
24 -----Including PCA.-----
25 best mean cv score= 0.6768973816233019
26 best parameters= {'metric': 'cosine', 'n_neighbors': 1}
27 [[530 17]
28  [ 13 17]]
29 0.767794028031688
30 ---This information is for the 4 th label.---
31 -----Including PCA.-----
32 best mean cv score= 0.639949171815369
33 best parameters= {'metric': 'euclidean', 'n_neighbors': 3}
34 [[512 22]
35  [ 58 21]]
36 0.6123121414687337
37 ---This information is for the 5 th label.---
38 -----Including PCA.-----
39 best mean cv score= 0.7037870122495336
40 best parameters= {'metric': 'euclidean', 'n_neighbors': 1}
41 [[636 18]
42  [ 20 13]]
43 0.6832082290797887
44 ---This information is for the 6 th label.---
45 -----Including PCA.-----
46 best mean cv score= 0.5970295738998852
47 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}

```

```

48 [[613   8]
49 [ 11   6]]
50 0.670029364402766
51 ---This information is for the 7 th label.---
52 -----Including PCA.-----
53 best mean cv score= 0.6628518308101344
54 best parameters= {'metric': 'cosine', 'n_neighbors': 1}
55 [[431  55]
56 [ 45  46]]
57 0.696162890607335
58 ---This information is for the 8 th label.---
59 -----Including PCA.-----
60 best mean cv score= 0.6327067238634185
61 best parameters= {'metric': 'euclidean', 'n_neighbors': 1}
62 [[660  15]
63 [ 17   8]]
64 0.6488888888888888
65 ---This information is for the 9 th label.---
66 -----Including PCA.-----
67 best mean cv score= 0.6296254938216598
68 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
69 [[578  26]
70 [ 22  13]]
71 0.6641911069063386
72 ---This information is for the 10 th label.---
73 -----Including PCA.-----
74 best mean cv score= 0.7639088766366591
75 best parameters= {'metric': 'cosine', 'n_neighbors': 1}
76 [[443  41]
77 [ 32  58]]
78 0.7798668503213958
79 ---This information is for the 11 th label.---
80 -----Including PCA.-----
81 best mean cv score= 0.6857203682733718
82 best parameters= {'metric': 'cosine', 'n_neighbors': 1}
83 [[608  21]
84 [ 23  18]]
85 0.7028190313699639

```

Listing 5: GridSearchCV-KNN-PCA results

I implemented the KNN algorithm without reducing the PCA dimension on the data, but the desired result is usually better on the KNN algorithm with the PCA dimension reduction algorithm. Although PCA improves the result in some algorithms such as KNN, on the other hand, some algorithms such as The Decision Tree algorithm can worsen the results. The reason for this is almost apparent given the nature of these two algorithms. Anyway, this is just my guess, and it does not mean that it is true. So I did not use PCA for the algorithms that I guessed PCA would not help.

In Listing 6, you can see the code related to running the KNN algorithm on the data without running the PCA algorithm. You can also compare the results of this performance in Listing 7 with the written results in the Listing 5.

```

1 print("KNeighborsClassifier:")
2 for i in np.arange(0,12,1):
3     knn = KNeighborsClassifier()
4     param_grid={'n_neighbors':np.arange(1,11,2),'metric':['euclidean','cosine','manhattan']}
5     grid=GridSearchCV(knn,param_grid=param_grid, cv=10,
6     scoring='balanced_accuracy', return_train_score=False)
7     grid.fit(X_training_data[i].to_numpy(),
8     y_training_data[i])
9     print("---This information is for the ",i,"th label.---")
10    print("best mean cv score= ",grid.best_score_)
11    print("best parameters= ",grid.best_params_)
12    y_pred_knn=grid.best_estimator_.predict(X_test[i])
13    print(confusion_matrix(y_test[i],np.transpose(np.
14        matrix(y_pred_knn))))

```

```

12     print(grid.score(X_test[i],y_test[i]))

```

Listing 6: GridSearchCV-KNN [5]

```

1 KNeighborsClassifier:
2 ---This information is for the 0 th label.---
3 best mean cv score= 0.7309100650238083
4 best parameters= {'metric': 'manhattan', 'n_neighbors': 9}
5 [[683  3]
6 [ 15  16]]
7 0.7558779272077494
8 ---This information is for the 1 th label.---
9 best mean cv score= 0.7846450333353117
10 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
11 [[634  12]
12 [ 6  17]]
13 0.8602772916947099
14 ---This information is for the 2 th label.---
15 best mean cv score= 0.7272648019793841
16 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
17 [[545  26]
18 [ 32  45]]
19 0.7694407169013123
20 ---This information is for the 3 th label.---
21 best mean cv score= 0.6781417541315211
22 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
23 [[528  19]
24 [ 13  17]]
25 0.7659658744667885
26 ---This information is for the 4 th label.---
27 best mean cv score= 0.6430757718521802
28 best parameters= {'metric': 'manhattan', 'n_neighbors': 3}
29 [[505  29]
30 [ 59  20]]
31 0.5994287204285782
32 ---This information is for the 5 th label.---
33 best mean cv score= 0.7049780678881433
34 best parameters= {'metric': 'euclidean', 'n_neighbors': 1}
35 [[637  17]
36 [ 20  13]]
37 0.6839727550736725
38 ---This information is for the 6 th label.---
39 best mean cv score= 0.6020852568911459
40 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
41 [[611  10]
42 [ 11  6]]
43 0.6684190584446339
44 ---This information is for the 7 th label.---
45 best mean cv score= 0.6631432384485604
46 best parameters= {'metric': 'cosine', 'n_neighbors': 1}
47 [[434  52]
48 [ 45  46]]
49 0.6992493103604215
50 ---This information is for the 8 th label.---
51 best mean cv score= 0.6377311937912507
52 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
53 [[660  15]
54 [ 19  6]]
55 0.6088888888888888
56 ---This information is for the 9 th label.---
57 best mean cv score= 0.6371383700958552
58 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
59 [[579  25]
60 [ 23  12]]
61 0.6507332071901608
62 ---This information is for the 10 th label.---
63 best mean cv score= 0.7758465460133529
64 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
65 [[447  37]
66 [ 31  59]]
67 0.78955463728191
68 ---This information is for the 11 th label.---
69 best mean cv score= 0.6895577029999881
70 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
71 [[610  19]
72 [ 22  19]]
```

73 0.7166039784404203

Listing 7: GridSearchCV-KNN results

## 6.2 Decision Tree

Decision Tree learning or induction of Decision Trees is one of the predictive modelling approaches used in statistics, data mining and machine learning. It uses a Decision Tree (as a predictive model) to go from observations about an item (represented in the branches) to conclusions about the item's target value (represented in the leaves). Tree models where the target variable can take a discrete set of values are called classification trees; in these tree structures, leaves represent class labels and branches represent conjunctions of features that lead to those class labels. Decision Trees where the target variable can take continuous values (typically real numbers) are called regression trees. Decision Trees are among the most popular machine learning algorithms given their intelligibility and simplicity. [10]

One of the algorithms I applied to the data (without PCA interference) was the Decision Tree algorithm. It should be noted that, as I mentioned earlier, my results show that PCA dimensionally reduced data do not make outstanding results with the Decision Tree. I think it's because the Decision Tree itself has a feature selection, but all the features are somehow important when we use PCA. However, because the results were not desirable, I decided to avoid writing the results of the Decision Tree algorithm on the data with PCA.

```

1 print("DecisionTreeClassifier:")
2 from sklearn.tree import DecisionTreeClassifier,
3     export_graphviz
4 for i in np.arange(0,12,1):
5     DecisionTree = DecisionTreeClassifier()
6     param_grid={'criterion':['gini','entropy'],
7                 'max_depth':np.arange(4,30,1)}
8     grid=GridSearchCV(DecisionTree,param_grid=
9                         param_grid, cv=10, scoring='balanced_accuracy',
10                        return_train_score=False)
11    grid.fit(X_training_data[i].to_numpy(),
12              y_training_data[i])
13    print("----This information is for the ",i,"th
14      label-----")
15    print("best mean cv score= ",grid.best_score_)
16    print("best parameters= ",grid.best_params_)
17    y_pred_dt=grid.best_estimator_.predict(X_test[i])
18    print(confusion_matrix(y_test[i],np.transpose(np.
19      matrix(y_pred_dt))))
20    print(grid.score(X_test[i],y_test[i]))
```

Listing 8: GridSearchCV-Decision Tree [5]

```

1 DecisionTreeClassifier:
2 -----This information is for the 0 th label
3 best mean cv score= 0.7343343187519519
4 best parameters= {'criterion': 'entropy', 'max_depth': 4}
5 [[682  4]
6  [ 15 16]]
7 0.7551490642339885
8 -----This information is for the 1 th label
9 best mean cv score= 0.7493516924842225
10 best parameters= {'criterion': 'gini', 'max_depth': 27}
11 [[637  9]
12  [ 7 16]]
13 0.8408601426840758
14 -----This information is for the 2 th label
15 best mean cv score= 0.715454431723965
16 best parameters= {'criterion': 'gini', 'max_depth': 27}
17 [[534  37]]
```

```

18 [[ 34 43]]
19 0.7468214797461732
20 -----This information is for the 3 th label
21 best mean cv score= 0.6357916597143412
22 best parameters= {'criterion': 'gini', 'max_depth': 22}
23 [[523 24]
24  [ 19 11]]
25 0.6613954905545399
26 -----This information is for the 4 th label
27 best mean cv score= 0.6307871144050721
28 best parameters= {'criterion': 'entropy', 'max_depth': 13}
29 [[487 47]
30  [ 53 26]]
31 0.6205494713886124
32 -----This information is for the 5 th label
33 best mean cv score= 0.6988894876519602
34 best parameters= {'criterion': 'entropy', 'max_depth': 28}
35 [[627 27]
36  [ 23 10]]
37 0.6308729496802892
38 -----This information is for the 6 th label
39 best mean cv score= 0.5994250886439558
40 best parameters= {'criterion': 'gini', 'max_depth': 28}
41 [[607 14]
42  [ 10 7]]
43 0.6946102112342522
44 -----This information is for the 7 th label
45 best mean cv score= 0.6527794346669296
46 best parameters= {'criterion': 'gini', 'max_depth': 21}
47 [[431 55]
48  [ 53 38]]
49 0.6522068466512911
50 -----This information is for the 8 th label
51 best mean cv score= 0.617416934552133
52 best parameters= {'criterion': 'gini', 'max_depth': 28}
53 [[652 23]
54  [ 18 7]]
55 0.6229629629629629
56 -----This information is for the 9 th label
57 best mean cv score= 0.6111641916613375
58 best parameters= {'criterion': 'entropy', 'max_depth': 21}
59 [[575 29]
60  [ 24 11]]
61 0.6331362346263009
62 -----This information is for the 10 th
63 label.
64 best mean cv score= 0.7517233591876147
65 best parameters= {'criterion': 'gini', 'max_depth': 20}
66 [[448 36]
67  [ 37 53]]
68 0.7572543617998164
69 -----This information is for the 11 th
70 label.
71 best mean cv score= 0.6584764787503303
72 best parameters= {'criterion': 'gini', 'max_depth': 28}
73 [[589 40]
74  [ 28 13]]
75 0.6267400829811159
```

Listing 9: GridSearchCV-Decision Tree results

## 6.3 Random forest

A Random forest is a classifier consisting of a collection of Decision Trees, where each tree is constructed by applying an algorithm A on the training set S and an additional random vector,  $\theta$ , where  $\theta$  is sampled i.i.d.<sup>8</sup> from some distribution. The prediction of the Random forest is obtained by a majority vote over the predictions of the individual trees.

With the argument I made for the Decision Tree algorithm, I will omit the Random forest algorithm and the rest of the algorithms written based on the Decision Tree on PCA data. Of course, I checked the background results and just ignored them for reporting. Finally, I would like to draw your attention to the implementation of the Random forest algorithm on the data in Listings 10 and 11.

```

1 print("RandomForestClassifier:")
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.datasets import make_classification
4 for i in np.arange(0,12,1):
5     RFC = RandomForestClassifier()
6     param_grid=[{'n_estimators':[500,700], 'criterion':['gini','entropy'],'max_depth':[4,6,8,10,15,20], 'max_features':['auto','sqrt','log2']}
7     grid=GridSearchCV(RFC,param_grid=param_grid, cv=5, scoring='balanced_accuracy', return_train_score=False)
8     grid.fit(X_training_data[i].to_numpy(), y_training_data[i])
9     print("-----This information is for the ",i,"th label-----")
10    print("best mean cv score= ",grid.best_score_)
11    print("best parameters= ",grid.best_params_)
12    y_pred_RFC=grid.best_estimator_.predict(X_test[i])
13    print(confusion_matrix(y_test[i],np.transpose(np.matrix(y_pred_RFC))))
14    print(grid.score(X_test[i],y_test[i]))

```

Listing 10: GridSearchCV-Random forest [5]

```

1 RandomForestClassifier:
2 -----This information is for the 0 th label
3 best mean cv score= 0.7216258089347315
4 best parameters= {'criterion': 'gini', 'max_depth': 15, 'max_features': 'auto', 'n_estimators': 700}
5 [[684  2]
6  [ 15 16]]
7 0.7566067901815103
8 -----This information is for the 1 th label
9 best mean cv score= 0.739999752677349
10 best parameters= {'criterion': 'entropy', 'max_depth': 10, 'max_features': 'sqrt', 'n_estimators': 500}
11 [[640   6]
12  [  6 17]]
13 0.8649212545430072
14 -----This information is for the 2 th label
15 best mean cv score= 0.6722003183198237
16 best parameters= {'criterion': 'entropy', 'max_depth': 20, 'max_features': 'auto', 'n_estimators': 700}
17 [[560  11]
18  [ 46 31]]
19 0.6916664771305752
20 -----This information is for the 3 th label
21 best mean cv score= 0.5665454586815841
22 best parameters= {'criterion': 'gini', 'max_depth': 20, 'max_features': 'sqrt', 'n_estimators': 500}
23 [[544   3]
24  [ 20 10]]
25 0.6639244363193175
26 -----This information is for the 4 th label
27 best mean cv score= 0.6182522327025808
28 best parameters= {'criterion': 'entropy', 'max_depth': 20, 'max_features': 'auto', 'n_estimators': 500}
29 [[526   8]
30  [ 57 22]]
31 0.631749869624994

```

```

32 -----This information is for the 5 th label
33 best mean cv score= 0.6767173252279635
34 best parameters= {'criterion': 'entropy', 'max_depth': 15, 'max_features': 'auto', 'n_estimators': 500}
35 [[652  2]
36  [ 23 10]]
37 0.649986099527384
38 -----This information is for the 6 th label
39 best mean cv score= 0.509426763957477
40 best parameters= {'criterion': 'gini', 'max_depth': 20, 'max_features': 'auto', 'n_estimators': 500}
41 [[621  0]
42  [ 17  0]]
43 0.5
44 -----This information is for the 7 th label
45 best mean cv score= 0.5917580102590566
46 best parameters= {'criterion': 'gini', 'max_depth': 20, 'max_features': 'sqrt', 'n_estimators': 500}
47 [[480   6]
48  [ 61 30]]
49 0.6586623253289919
50 -----This information is for the 8 th label
51 best mean cv score= 0.5349558244652022
52 best parameters= {'criterion': 'entropy', 'max_depth': 20, 'max_features': 'auto', 'n_estimators': 500}
53 [[674   1]
54  [ 23  2]]
55 0.5392592592592592
56 -----This information is for the 9 th label
57 best mean cv score= 0.5541034155597723
58 best parameters= {'criterion': 'gini', 'max_depth': 20, 'max_features': 'auto', 'n_estimators': 500}
59 [[603   1]
60  [ 30  5]]
61 0.570600756859035
62 -----This information is for the 10 th label
63 best mean cv score= 0.7399732819343566
64 best parameters= {'criterion': 'gini', 'max_depth': 20, 'max_features': 'auto', 'n_estimators': 500}
65 [[469  15]
66  [ 41 49]]
67 0.7567263544536271
68 -----This information is for the 11 th label
69 best mean cv score= 0.5619639353507783
70 best parameters= {'criterion': 'entropy', 'max_depth': 20, 'max_features': 'sqrt', 'n_estimators': 500}
71 [[628   1]
72  [ 35  6]]
73 0.5723758191476986

```

Listing 11: GridSearchCV-Random forest results

## 6.4 XGBoost

XGBoost is an optimized distributed gradient boosting library designed to be highly efficient, flexible and portable. It implements machine learning algorithms under the Gradient Boosting framework. XGBoost provides a parallel tree boosting (also known as GBDT, GBM) that solve many data science problems in a fast and accurate way. The same code runs on major distributed environment (Hadoop, SGE, MPI) and can solve problems beyond billions of examples. [11] In the following, I have implemented this beautiful algorithm on my data.

```

1 from xgboost import XGBClassifier
2 print("XGBClassifier:")
3 for i in np.arange(0,12,1):

```

```

4 param_grid={'n_estimators':[700,2000], 'max_depth'
5   : [3,4,5,6,7], 'eval_metric':['mlogloss']}
6 grid=GridSearchCV(XGBClassifier(),param_grid=
7 param_grid, cv=10,scoring='balanced_accuracy',
8 return_train_score=False)
9 grid.fit(X_training_data[i].to_numpy(),
10 y_training_data[i])
11 print("----This information is for the ",i,"th
12 label.-----")
13 print("best mean cv score= ",grid.best_score_)
14 print("best parameters= ",grid.best_params_)
15 y_pred_xgbc=grid.best_estimator_.predict(X_test[i])
16 print(confusion_matrix(y_test[i],np.transpose(np.
17 matrix(y_pred_xgbc))))
18 print(grid.score(X_test[i],y_test[i]))

```

Listing 12: GridSearchCV-XGBoost [5]

```

49 0.7115723782390448
50 ---This information is for the 8 th label.-----
51 best mean cv score= 0.597032747111124
52 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 4,
53   'n_estimators': 700}
54 [[669 6]
55 [ 21 4]]
56 0.5755555555555556
57 ---This information is for the 9 th label.-----
58 best mean cv score= 0.5932812079563602
59 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 6,
60   'n_estimators': 700}
61 [[601 3]
62 [ 27 8]]
63 0.611802270577105
64 ---This information is for the 10 th label.-----
65 best mean cv score= 0.7958197053234901
66 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 3,
67   'n_estimators': 700}
68 [[462 22]
69 [ 33 57]]
70 0.793939393939394
71 ---This information is for the 11 th label.-----
72 best mean cv score= 0.619645691038506
73 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 3,
74   'n_estimators': 700}
75 [[624 5]
76 [ 30 11]]
77 0.6301717786653225

```

Listing 13: GridSearchCV-XGBoost results

```

1 XGBClassifier:
2 ---This information is for the 0 th label.-----
3 best mean cv score= 0.7220923302944042
4 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 3,
5   'n_estimators': 700}
6 [[683 3]
7 [ 16 15]]
8 0.7397488949496849
9 ---This information is for the 1 th label.-----
10 best mean cv score= 0.7588862345935466
11 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 6,
12   'n_estimators': 2000}
13 [[640 6]
14 [ 8 15]]
15 0.821442993673442
16 ---This information is for the 2 th label.-----
17 best mean cv score= 0.7333850964654818
18 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 5,
19   'n_estimators': 2000}
20 [[557 14]
21 [ 36 41]]
22 0.7539745718379693
23 ---This information is for the 3 th label.-----
24 best mean cv score= 0.6181325689242583
25 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 5,
26   'n_estimators': 700}
27 [[540 7]
28 [ 20 10]]
29 0.6602681291895186
30 ---This information is for the 4 th label.-----
31 best mean cv score= 0.6417879130291102
32 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 3,
33   'n_estimators': 2000}
34 [[519 15]
35 [ 54 25]]
36 0.6441829042810412
37 ---This information is for the 5 th label.-----
38 best mean cv score= 0.6939448539211255
39 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 7,
40   'n_estimators': 700}
41 [[650 4]
42 [ 22 11]]
43 0.6636085626911314
44 ---This information is for the 6 th label.-----
45 best mean cv score= 0.5692212315899486
46 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 3,
47   'n_estimators': 700}
48 [[621 0]
49 [ 16 1]]
50 0.5294117647058824
51 ---This information is for the 7 th label.-----
52 best mean cv score= 0.6742232870501532
53 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 5,
54   'n_estimators': 2000}
55 [[462 24]
56 [ 48 43]]
```

## 6.5 C-Support Vector Classification

In machine learning, support-vector machines (SVMs, also support-vector networks[12]) are supervised learning models with associated learning algorithms that analyze data for classification and regression analysis. Developed at ATT Bell Laboratories by Vladimir Vapnik with colleagues.

```

1 from sklearn.svm import SVC
2 print("SVC:")
3 for i in np.arange(0,12,1):
4   param_grid={'kernel':['poly'],'degree':np.arange
5     (1,15,1)}
6   grid=GridSearchCV(SVC(),param_grid=param_grid,cv
7     =10,scoring='balanced_accuracy',return_train_score
8     =False)
9   grid.fit(X_training_data[i].to_numpy(),
10   y_training_data[i])
11   print("----This information is for the ",i,"th
12   label.-----")
13   print("best mean cv score= ",grid.best_score_)
14   print("best parameters= ",grid.best_params_)
15   y_pred_svm=grid.best_estimator_.predict(X_test[i])
16   print(confusion_matrix(y_test[i],np.transpose(np.
17   matrix(y_pred_svm))))
18   print(grid.score(X_test[i],y_test[i]))

```

Listing 14: GridSearchCV-SVC [5]

```

1 SVC:
2 ---This information is for the 0 th label.-----
3 best mean cv score= 0.7230681963217643
4 best parameters= {'degree': 4, 'kernel': 'poly'}
5 [[680 6]
6 [ 15 16]]
7 0.7536913382864667
8 ---This information is for the 1 th label.-----
9 best mean cv score= 0.7546105110095553
10 best parameters= {'degree': 5, 'kernel': 'poly'}
11 [[638 8]
12 [ 8 15]]
13 0.8198950060573429

```

```

14 ---This information is for the 2 th label.-----
15 best mean cv score= 0.7242156849858815
16 best parameters= {'degree': 5, 'kernel': 'poly'}
17 [[546 25]
18 [ 38 39]]
19 0.7313553346828303
20 ---This information is for the 3 th label.-----
21 best mean cv score= 0.6645036425232391
22 best parameters= {'degree': 7, 'kernel': 'poly'}
23 [[534 13]
24 [ 17 13]]
25 0.7047836684948202
26 ---This information is for the 4 th label.-----
27 best mean cv score= 0.630630697298831
28 best parameters= {'degree': 11, 'kernel': 'poly'}
29 [[490 44]
30 [ 56 23]]
31 0.6043711183805054
32 ---This information is for the 5 th label.-----
33 best mean cv score= 0.6947422904425824
34 best parameters= {'degree': 8, 'kernel': 'poly'}
35 [[645 9]
36 [ 20 13]]
37 0.6900889630247429
38 ---This information is for the 6 th label.-----
39 best mean cv score= 0.5878664257410506
40 best parameters= {'degree': 9, 'kernel': 'poly'}
41 [[612 9]
42 [ 11 6]]
43 0.6692242114236999
44 ---This information is for the 7 th label.-----
45 best mean cv score= 0.6656936034198049
46 best parameters= {'degree': 14, 'kernel': 'poly'}
47 [[433 53]
48 [ 54 37]]
49 0.6487699543255099
50 ---This information is for the 8 th label.-----
51 best mean cv score= 0.5954265276466763
52 best parameters= {'degree': 10, 'kernel': 'poly'}
53 [[667 8]
54 [ 18 7]]
55 0.634074074074074
56 ---This information is for the 9 th label.-----
57 best mean cv score= 0.6096192244051426
58 best parameters= {'degree': 10, 'kernel': 'poly'}
59 [[583 21]
60 [ 24 11]]
61 0.6397587511825922
62 ---This information is for the 10 th label.-----
63 best mean cv score= 0.7740261034794257
64 best parameters= {'degree': 5, 'kernel': 'poly'}
65 [[454 30]
66 [ 36 54]]
67 0.7690082644628099
68 ---This information is for the 11 th label.-----
69 best mean cv score= 0.6622435935009551
70 best parameters= {'degree': 7, 'kernel': 'poly'}
71 [[618 11]
72 [ 24 17]]
73 0.6985730350149288

```

Listing 15: GridSearchCV-SVC results

## 6.6 AdaBoost

AdaBoost, short for Adaptive Boosting, is a statistical classification meta-algorithm formulated by Yoav Freund and Robert Schapire, who won the 2003 Gödel Prize for their work. It can be used in conjunction with many other types of learning algorithms to improve performance. The output of the other learning algorithms ('weak learners') is combined into a weighted sum that represents the final output of the boosted classifier. AdaBoost is adaptive in the sense that subsequent weak learners are tweaked in favor of those instances misclassified by previous classifiers. In some problems it can be less susceptible to the overfitting

problem than other learning algorithms. The individual learners can be weak, but as long as the performance of each one is slightly better than random guessing, the final model can be proven to converge to a strong learner. [13]

```

1 from sklearn.ensemble import AdaBoostClassifier
2 print("AdaBoostClassifier:")
3 for i in np.arange(0,12,1):
4     param_grid={'n_estimators':[8000,1000,500,700],
5                 'algorithm':['SAMME.R','SAMME'],'learning_rate':
6                 :[1]}
7     grid=GridSearchCV(AdaBoostClassifier(),param_grid=
8                         param_grid, cv=5,scoring='balanced_accuracy',
9                         return_train_score=False)
10    grid.fit(X_training_data[i].to_numpy(),
11              y_training_data[i])
12    print("---This information is for the ",i,"th
13         label.-----")
14    print("best mean cv score= ",grid.best_score_)
15    print("best parameters= ",grid.best_params_)
16    y_pred_AdaBoost=grid.best_estimator_.predict(
17        X_test[i])
18    print(confusion_matrix(y_test[i],np.transpose(np.
19        matrix(y_pred_AdaBoost))))
20    print(grid.score(X_test[i],y_test[i]))

```

Listing 16: GridSearchCV-AdaBoost [5]

```

1 AdaBoostClassifier:
2 ---This information is for the 0 th label.-----
3 best mean cv score= 0.7157613477540118
4 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
5                   'n_estimators': 8000}
6 [[674 12]
7 [ 16 15]]
8 0.7331891281858365
9 ---This information is for the 1 th label.-----
10 best mean cv score= 0.7341166085078775
11 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
12                   'n_estimators': 500}
13 [[640 6]
14 [ 6 17]]
15 0.8649212545430072
16 ---This information is for the 2 th label.-----
17 best mean cv score= 0.7308967630172937
18 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
19                   'n_estimators': 8000}
20 [[535 36]
21 [ 35 42]]
22 0.7412036299952236
23 ---This information is for the 3 th label.-----
24 best mean cv score= 0.6154114251647196
25 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
26                   'n_estimators': 8000}
27 [[534 13]
28 [ 16 14]]
29 0.7214503351614869
30 ---This information is for the 4 th label.-----
31 best mean cv score= 0.6233451764286004
32 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
33                   'n_estimators': 1000}
34 [[511 23]
35 [ 57 22]]
36 0.6177049258047693
37 ---This information is for the 5 th label.-----
38 best mean cv score= 0.6699627297727602
39 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
40                   'n_estimators': 8000}
41 [[643 11]
42 [ 21 12]]
43 0.6734083958854601
44 ---This information is for the 6 th label.-----
45 best mean cv score= 0.5620798156865388
46 best parameters= {'algorithm': 'SAMME', 'learning_rate': 1,
47                   'n_estimators': 8000}

```

```

41 [[617  4]                                     ---This information is for the 0 th label.-----
42 [ 17  0]]                                     best mean cv score= 0.7269508786019065
43 0.4967793880837359
44 ---This information is for the 7 th label.-----
45 best mean cv score= 0.6537217100575029
46 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
   'n_estimators': 8000}
47 [[432  54]                                     ---This information is for the 1 th label.-----
48 [ 54 37]]                                     best mean cv score= 0.7326826795066174
49 0.6477411477411477
50 ---This information is for the 8 th label.-----
51 best mean cv score= 0.5969505584559676
52 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
   'n_estimators': 8000}
53 [[658  17]                                     ---This information is for the 2 th label.-----
54 [ 18  7]]                                     best mean cv score= 0.6652733608745625
55 0.6274074074074074
56 ---This information is for the 9 th label.-----
57 best mean cv score= 0.5858266150853018
58 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
   'n_estimators': 700}
59 [[600  4]                                      ---This information is for the 3 th label.-----
60 [ 29  6]]                                     best mean cv score= 0.559236688270835
61 0.58240302743614
62 ---This information is for the 10 th label.-----
63 best mean cv score= 0.7804093635981373
64 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
   'n_estimators': 1000}
65 [[458  26]                                     ---This information is for the 4 th label.-----
66 [ 36 54]]                                     best mean cv score= 0.6169766928272998
67 0.7731404958677686
68 ---This information is for the 11 th label.-----
69 best mean cv score= 0.624099446415163
70 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
   'n_estimators': 1000}
71 [[617  12]                                     ---This information is for the 5 th label.-----
72 [ 32  9]]                                     best mean cv score= 0.6717171804892169
73 0.6002171468455543

```

Listing 17: GridSearchCV-AdaBoost results

## 6.7 Bagging

Bootstrap aggregating, also called bagging (from bootstrap aggregating), is a machine learning ensemble meta-algorithm designed to improve the stability and accuracy of machine learning algorithms used in statistical classification and regression. It also reduces variance and helps to avoid overfitting<sup>9</sup>. Although it is usually applied to Decision Tree methods, it can be used with any type of method. Bagging is a special case of the model averaging approach. [14]

```

1 from sklearn.ensemble import BaggingClassifier
2 print("BaggingClassifier:")
3 for i in np.arange(0,12,1):
4     param_grid={'n_estimators':[10,100,500],
5      'max_samples':[1,0.5],'max_features':[1,0.8,0.6]}
6     grid=GridSearchCV(BaggingClassifier(),param_grid=
7     param_grid, cv=5,scoring='balanced_accuracy',
8     return_train_score=False)
9     grid.fit(X_training_data[i].to_numpy(),
10    y_training_data[i])
11    print("----This information is for the ",i,"th
12    label.-----")
13    print("best mean cv score= ",grid.best_score_)
14    print("best parameters= ",grid.best_params_)
15    y_pred_bag=grid.best_estimator_.predict(X_test[i])
16    print(confusion_matrix(y_test[i],np.transpose(np.
17      matrix(y_pred_bag))))
18    print(grid.score(X_test[i],y_test[i]))

```

Listing 18: GridSearchCV-Bagging [5]

```

2 ---This information is for the 0 th label.-----
3 best mean cv score= 0.7269508786019065
4 best parameters= {'max_features': 0.6, 'max_samples': 0.5,
   'n_estimators': 500}
5 [[683  3]
6 [ 15 16]]                                     ---This information is for the 1 th label.-----
7 0.7558779272077494
8 best mean cv score= 0.7326826795066174
9 best parameters= {'max_features': 0.8, 'max_samples': 0.5,
   'n_estimators': 100}
10 [[640  6]
11 [ 8 15]]                                     ---This information is for the 2 th label.-----
12 0.821442993673442
13 ---This information is for the 3 th label.-----
14 best mean cv score= 0.6652733608745625
15 best parameters= {'max_features': 0.6, 'max_samples': 0.5,
   'n_estimators': 100}
16 [[560  11]
17 [ 48 29]]                                     ---This information is for the 4 th label.-----
18 0.6786794641435623
19 ---This information is for the 5 th label.-----
20 best mean cv score= 0.559236688270835
21 best parameters= {'max_features': 0.8, 'max_samples': 0.5,
   'n_estimators': 500}
22 [[545  2]
23 [ 24 6]]                                     ---This information is for the 6 th label.-----
24 0.5981718464351006
25 ---This information is for the 7 th label.-----
26 best mean cv score= 0.6169766928272998
27 best parameters= {'max_features': 0.6, 'max_samples': 0.5,
   'n_estimators': 100}
28 [[524  10]
29 [ 57 22]]                                     ---This information is for the 8 th label.-----
30 0.629877210448964
31 ---This information is for the 9 th label.-----
32 best mean cv score= 0.6717171804892169
33 best parameters= {'max_features': 0.8, 'max_samples': 0.5,
   'n_estimators': 100}
34 [[653  1]
35 [ 20 13]]                                     ---This information is for the 10 th label.-----
36 0.6962051709758132
37 ---This information is for the 11 th label.-----
38 best mean cv score= 0.5288716095018979
39 best parameters= {'max_features': 0.8, 'max_samples': 0.5,
   'n_estimators': 10}
40 [[621  0]
41 [ 17 0]]                                     ---This information is for the 12 th label.-----
42 0.5
43 ---This information is for the 13 th label.-----
44 best mean cv score= 0.5930241135390106
45 best parameters= {'max_features': 0.8, 'max_samples': 0.5,
   'n_estimators': 100}
46 [[478  8]
47 [ 62 29]]                                     ---This information is for the 14 th label.-----
48 0.6511102066657621
49 ---This information is for the 15 th label.-----
50 best mean cv score= 0.529686773652454
51 best parameters= {'max_features': 0.8, 'max_samples': 0.5,
   'n_estimators': 10}
52 [[674  1]
53 [ 24 1]]                                     ---This information is for the 16 th label.-----
54 0.5192592592592592
55 ---This information is for the 17 th label.-----
56 best mean cv score= 0.5476517181011051
57 best parameters= {'max_features': 0.8, 'max_samples': 0.5,
   'n_estimators': 500}
58 [[603  1]
59 [ 29 6]]                                     ---This information is for the 18 th label.-----
60 0.5848864711447493
61 ---This information is for the 19 th label.-----
62 best mean cv score= 0.7392036815039424
63 best parameters= {'max_features': 0.8, 'max_samples': 0.5,
   'n_estimators': 500}
64 [[466  18]
65 [ 42 48]]                                     ---This information is for the 20 th label.-----
66 0.7402036815039424

```

```

67 0.7480716253443527
68 ---This information is for the 11 th label.-----
69 best mean cv score= 0.5421860682323457
70 best parameters= {'max_features': 0.8, 'max_samples': 0.5, 'n_estimators': 100}
71 [[628 1]
72 [ 35 6]]
73 0.5723758191476986

```

Listing 19: GridSearchCV-Bagging results

## 6.8 MLP

A multilayer perceptron (MLP) is a class of feedforward artificial neural network (ANN). The term MLP is used ambiguously, sometimes loosely to mean any feedforward ANN, sometimes strictly to refer to networks composed of multiple layers of perceptrons (with threshold activation); see § Terminology. Multilayer perceptrons are sometimes colloquially referred to as "vanilla" neural networks, especially when they have a single hidden layer.[15]

An MLP consists of at least three layers of nodes: an input layer, a hidden layer and an output layer. Except for the input nodes, each node is a neuron that uses a nonlinear activation function. MLP utilizes a supervised learning technique called backpropagation for training.[16][17] Its multiple layers and non-linear activation distinguish MLP from a linear perceptron. It can distinguish data that is not linearly separable.[18]

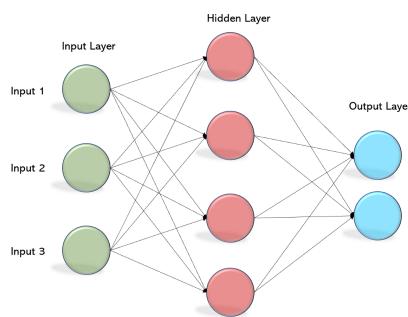


Figure 6: MLP neural network graph

```

1 from sklearn.neural_network import MLPClassifier
2 print("MLPClassifier:")
3 for i in np.arange(0,12,1):
4     param_grid={'activation':['identity', 'logistic',
5         'tanh', 'relu'],'hidden_layer_sizes':[100,500]}
6     grid=GridSearchCV(MLPClassifier(),param_grid=
7         param_grid, cv=5, scoring='balanced_accuracy',
8         return_train_score=False)
9     grid.fit(X_training_data[i].to_numpy(),
10         y_training_data[i])
11    print(" ---This information is for the ",i,"th
12    label.-----")
13    print("best mean cv score= ",grid.best_score_)
14    print("best parameters= ",grid.best_params_)
15    y_pred_bag=grid.best_estimator_.predict(X_test[i])
16    print(confusion_matrix(y_test[i],np.transpose(np.
17        matrix(y_pred_bag))))
18    print(grid.score(X_test[i],y_test[i]))

```

Listing 20: GridSearchCV-MLP [5]

```

1 MLPClassifier:
2 ---This information is for the 0 th label.-----
3 best mean cv score= 0.7190012640264282

```

```

4 best parameters= {'activation': 'logistic', 'hidden_layer_sizes': 500}
5 [[684 2]
6 [ 15 16]]
7 0.7566067901815103
8 ---This information is for the 1 th label.-----
9 best mean cv score= 0.7713046063735265
10 best parameters= {'activation': 'relu', 'hidden_layer_sizes': 100}
11 [[639 7]
12 [ 7 16]]
13 0.8424081303001749
14 ---This information is for the 2 th label.-----
15 best mean cv score= 0.7485781227406563
16 best parameters= {'activation': 'tanh', 'hidden_layer_sizes': 500}
17 [[544 27]
18 [ 36 41]]
19 0.7425910341847295
20 ---This information is for the 3 th label.-----
21 best mean cv score= 0.6675002192944698
22 best parameters= {'activation': 'relu', 'hidden_layer_sizes': 500}
23 [[540 7]
24 [ 18 12]]
25 0.6936014625228519
26 ---This information is for the 4 th label.-----
27 best mean cv score= 0.6513103184676161
28 best parameters= {'activation': 'tanh', 'hidden_layer_sizes': 500}
29 [[511 23]
30 [ 57 22]]
31 0.6177049258047693
32 ---This information is for the 5 th label.-----
33 best mean cv score= 0.6968680706325083
34 best parameters= {'activation': 'relu', 'hidden_layer_sizes': 100}
35 [[646 8]
36 [ 20 13]]
37 0.6908534890186266
38 ---This information is for the 6 th label.-----
39 best mean cv score= 0.5938001655740702
40 best parameters= {'activation': 'tanh', 'hidden_layer_sizes': 500}
41 [[617 4]
42 [ 13 4]]
43 0.6144264469072653
44 ---This information is for the 7 th label.-----
45 best mean cv score= 0.6845563149901454
46 best parameters= {'activation': 'tanh', 'hidden_layer_sizes': 100}
47 [[461 25]
48 [ 55 36]]
49 0.6720820331931443
50 ---This information is for the 8 th label.-----
51 best mean cv score= 0.6304795924052781
52 best parameters= {'activation': 'tanh', 'hidden_layer_sizes': 100}
53 [[663 12]
54 [ 19 6]]
55 0.6111111111111112
56 ---This information is for the 9 th label.-----
57 best mean cv score= 0.6378664632826452
58 best parameters= {'activation': 'tanh', 'hidden_layer_sizes': 500}
59 [[594 10]
60 [ 27 8]]
61 0.6060075685903501
62 ---This information is for the 10 th label.-----
63 best mean cv score= 0.798803902979208
64 best parameters= {'activation': 'tanh', 'hidden_layer_sizes': 100}
65 [[452 32]
66 [ 34 56]]
67 0.7780532598714417
68 ---This information is for the 11 th label.-----

```

```

69 best_mean_cv_score= 0.6682045890004455
70 best_parameters= {'activation': 'tanh', 'hidden_layer_sizes':
    500}
71 [[622 7]
72 [ 23 18]]
73 0.7139478072046221

```

Listing 21: GridSearchCV-MLP results

## 7 EVALUATION

In this section, I want to compare the accuracy of different algorithms for each label that we have obtained based on balanced accuracy. First of all, let's have a summary of what was done. I prepared the data to be processed by different algorithms. Then I ran different algorithms on the data. The results included cv score, confusion matrix, and the absolute accuracy of the algorithm on the test set. In this section, I will select the best predictor algorithm for each dimension of the label vector based on the results.

For the head dimension of the label vector, I prepared a figure containing three diagrams [such as Figure 7]. The first graph is the accuracy obtained from the test set based on balanced accuracy. The second diagram contains the cv score output of each algorithm.<sup>10</sup> The third diagram compares the previous two, giving us great information about overfitting or underfitting. In general, the most important parameter to choose the best algorithm will be balanced accuracy, then the analysis of the third diagram, cv score, and confusion matrix can help us in selecting the best choice.

### 7.1 1st dimension of the labels vector

Figure 7 shows us that almost all algorithms have been able to be accurate at practically the same level. However, the highest accuracy was on the test set for Random forest and multi layer perceptron algorithms. Fortunately, all test set results are better than cv scores, which means we do not have overfitting. The AdaBoost algorithm also had the lowest accuracy on the test set.

Which is better between Random forest and multi layer perceptron algorithms? Random forest seems to be a good option due to the use of fewer features and a higher cv score. But the execution runtime of the multilayer perceptron algorithm is significantly shorter. Assuming we do not have a computation time problem, I select the Random forest algorithm with the hyperparameters mentioned in Listing 11.

### 7.2 2nd dimension of the labels vector

Figure 8 also shows us that almost all algorithms have the same level of accuracy. However, the highest accuracy on the test set was for Random forest and AdaBoost algorithms. I'm glad that in the second dimension, as in the first dimension, all the results related to the test set are better than the cv scores, and this, fortunately, means that we are not overfitting. Also, the SVC algorithm had the lowest accuracy on the test set.

Which is better between Random forest and AdaBoost algorithms? The Random forest seems to be a good option due to the use of fewer features and a higher cv score. Also, the execution time of the AdaBoost algorithm is not much different from the Random forest. I chose the Random forest algorithm with the hyperparameters mentioned in Listing 11 as the best algorithm for estimating the second dimension of the label vector.

Finally, I should note that the selected algorithm has achieved 86.49212545430072 % balanced accuracy.

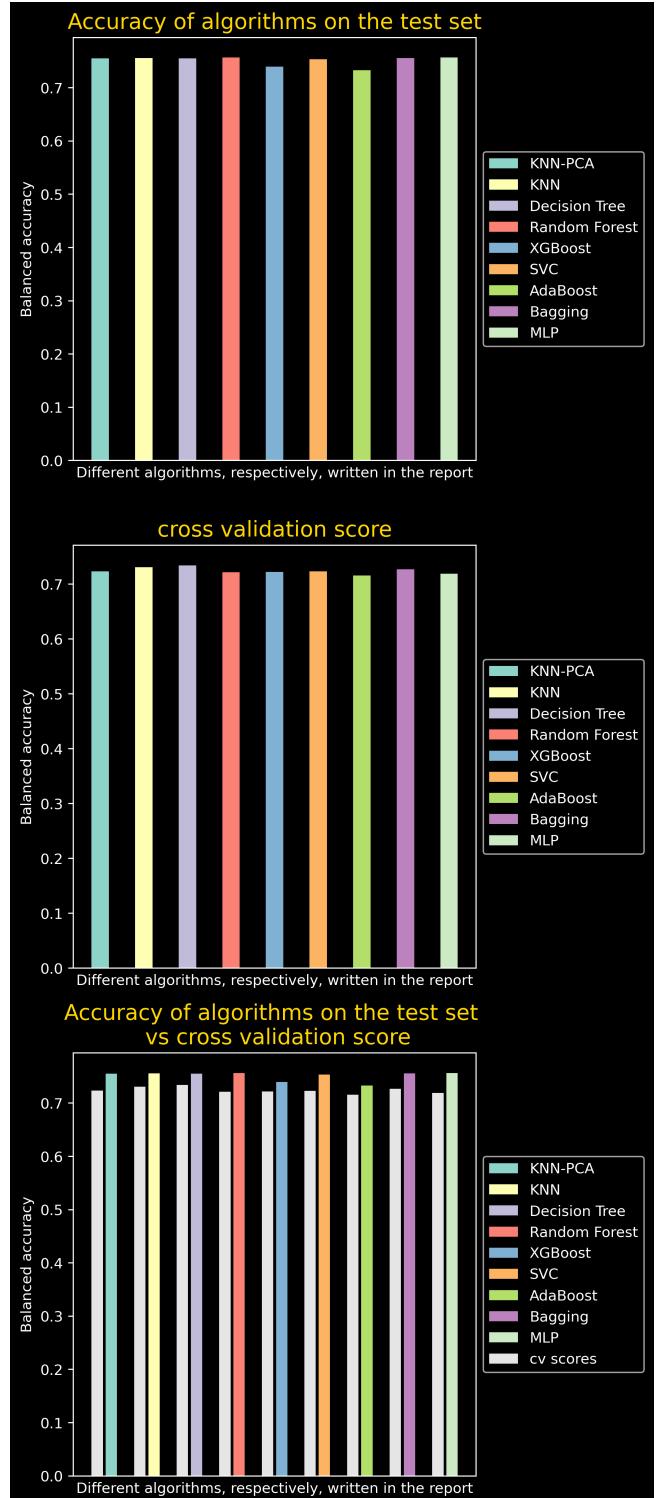


Figure 7: Diagrams of the 1st dimension of the labels vectors[4]

### 7.3 3rd dimension of the labels vector

Figure 9 shows that the accuracy of algorithms has more variance than the first and second dimensions. However, the highest accuracy has been on the test set for KNN algorithm. Except for the MLP algorithm, the test set results are better than the cv scores in the third dimension, which means we are not likely to see overfitting in this dimension. The Bagging algorithm also had the lowest accuracy on the test set.

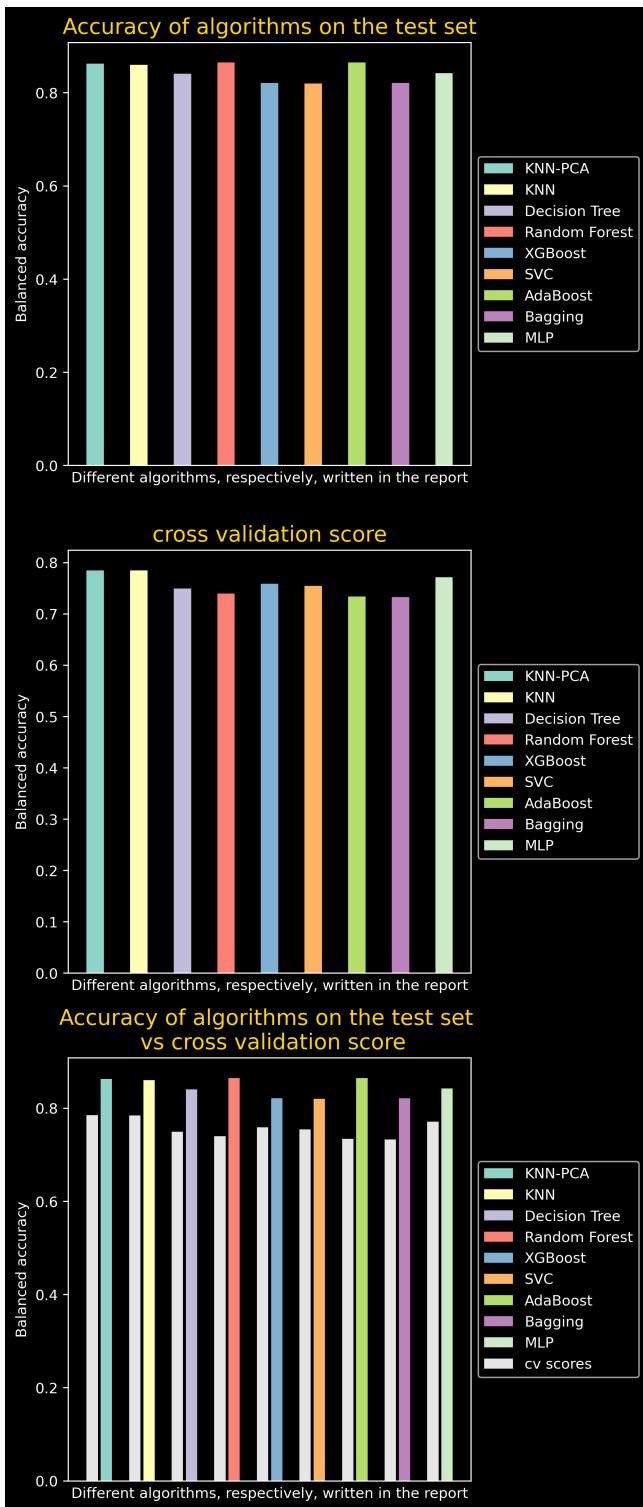


Figure 8: Diagrams of the 2nd dimension of the labels vectors[4]

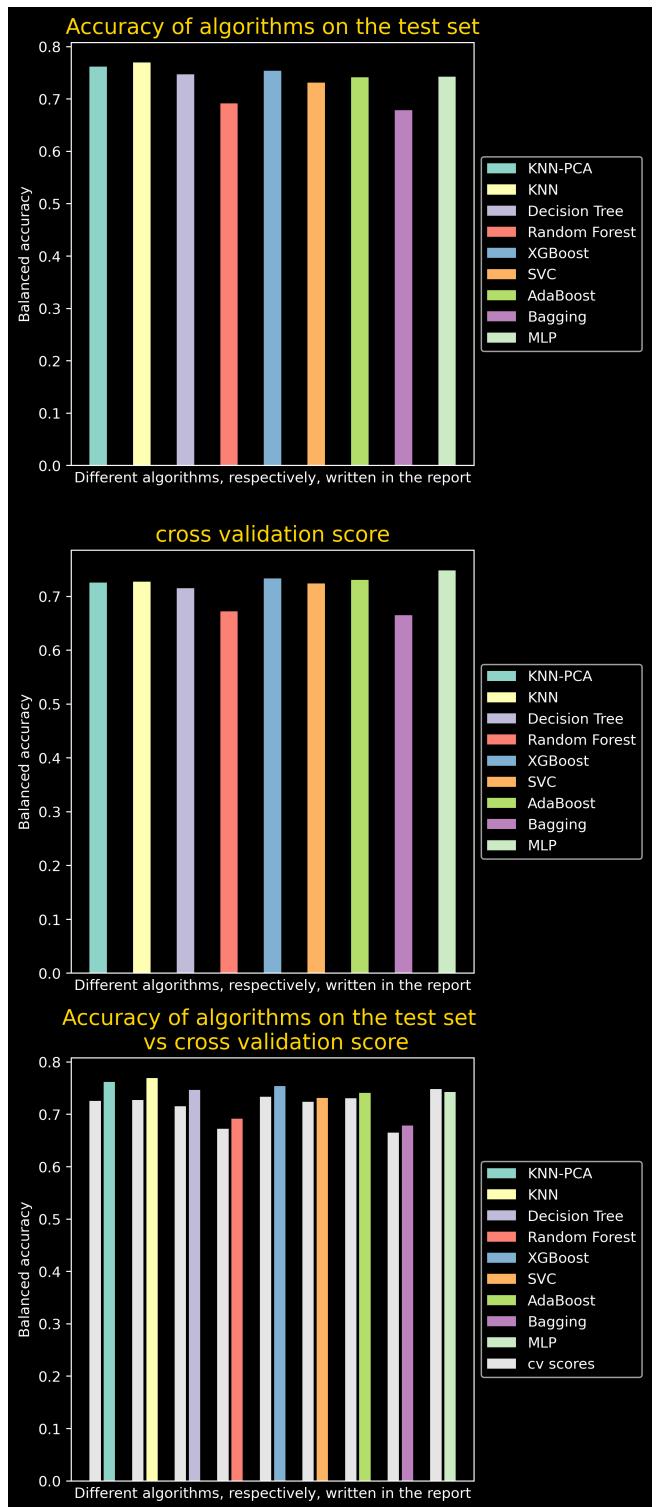


Figure 9: Diagrams of the 3rd dimension of the labels vectors[4]

#### 7.4 4th dimension of the labels vector

Interestingly, KNN has performed better on this data dimension than KNN-PCA. So I was curious to see how KNN performed on all dimensions of my data? I found that KNN did better in 6 dimensions and worse in 6 dimensions. But adding up their error difference, I found that overall KNN performance was slightly weaker than KNN-PCA. [See cv scores, which means we are not trapped in overfitting. The Bagging algorithm also had the lowest accuracy on the test set. Figure 10]

KNN-PCA works better than KNN in dimensions with grayscale colors.

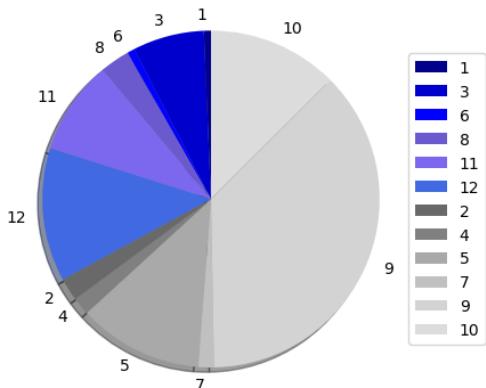


Figure 10: KNN vs. KNN-PCA performance in this project[4]

The interesting thing is that in the previous dimension where KNN performed well, Bagging still had the weakest performance.

## 7.5 5th dimension of the labels vector

Figure 12 shows almost a slight variance between the accuracy of the different algorithms. However, the most careful focus was on the test set for the XGBoost algorithm. Unfortunately, in the fifth dimension, we see some algorithms with a weak performance in test set score vs. cv scores. The KNN algorithm also had the lowest accuracy on the test set.

Of course, the highest cv score is related to the MLP algorithm, which, as you can see in Figure 12, had a poorer performance in the test set than the cv score. I think it smells overfitting. As a result, we leave this choice aside. In the dimension of the XGBoost algorithm achieved 64.41829042810412% balanced accuracy.

## 7.6 6th dimension of the labels vector

According to Figure 5, apart from the Bagging algorithm's highest test score in the sixth dimension, only the two Bagging and AdaBoost algorithms have obtained a higher test score than their cv score. This is while the average balanced accuracy of all algorithms is 67.357829%. And the AdaBoost algorithm has a slightly lower test score than the average. But the Bagging algorithm has a test score that is more than 2.2% higher than the average. For these reasons, it is enough to consider Bagging as the best predictor for sixth dimension labels.

## 7.7 7th dimension of the labels vector

Figure 14 shows the relative success of simpler algorithms compared to time-consuming algorithms. The seventh dimension is where the Decision Tree algorithm shows its power. Looking at Figure 22, you can see how the Decision Tree algorithm differs from other algorithms much better. The average balanced accuracy of all executed algorithms is more than 10% less than the accuracy obtained by the Decision Tree algorithm.

It is worth noting that even Decision Tree-based algorithms, such as Random forest and Bagging, have exactly the same confusion matrices. They could not even predict a True label number correctly.

## 7.8 8th dimension of the labels vector

Figure 15 shows the highest accuracy on the test set for XGBoost algorithms. In the eighth dimension, unfortunately, 33% of the algorithms in the test set were less accurate than the cv score. The

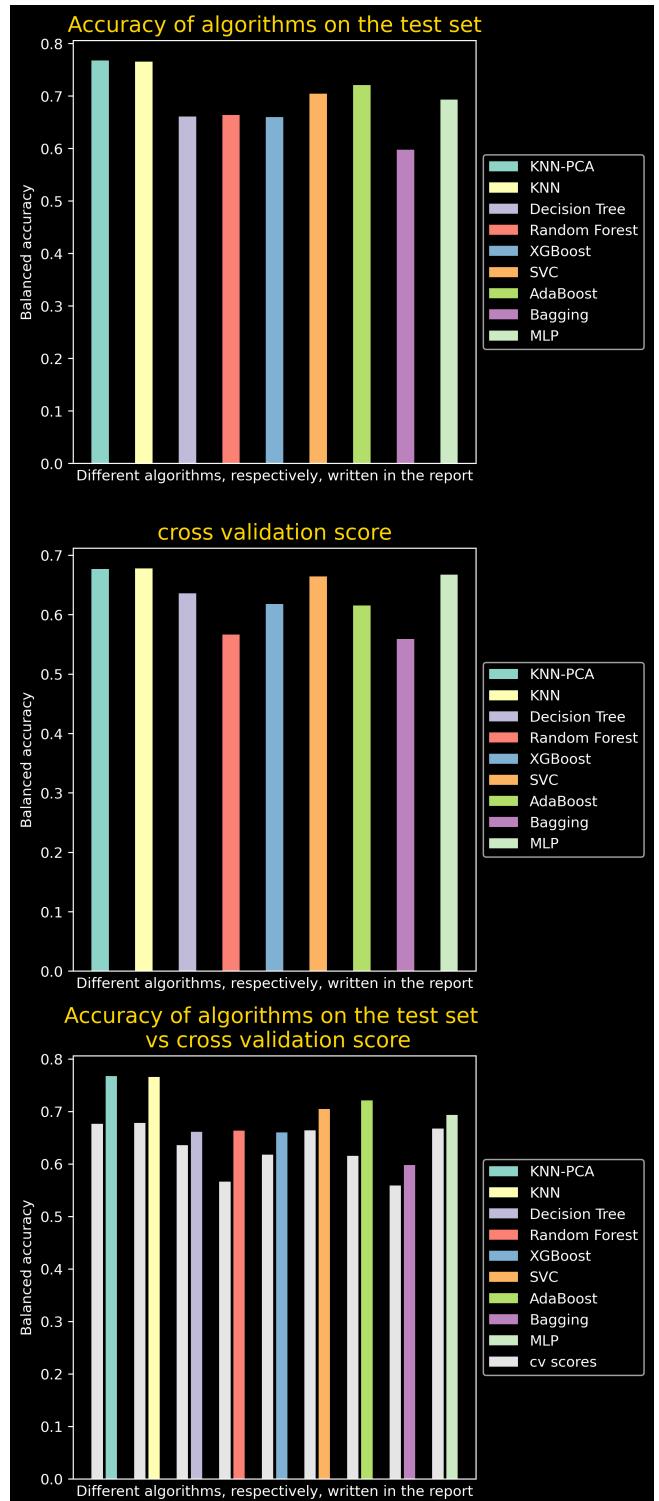


Figure 11: Diagrams of the 4th dimension of the labels vectors[4]

Bagging algorithm also had the lowest accuracy on the test set. The KNN-PCA and KNN algorithms performed well but were not the best. We have already seen the inverse relationship of the Bagging algorithm with these two algorithms.

It is not harmful to know that there is about a 99% correlation between the accuracy of KNN-PCA and KNN algorithms. Interestingly, the lowest correlation between the total accuracy of the implemented algorithms is

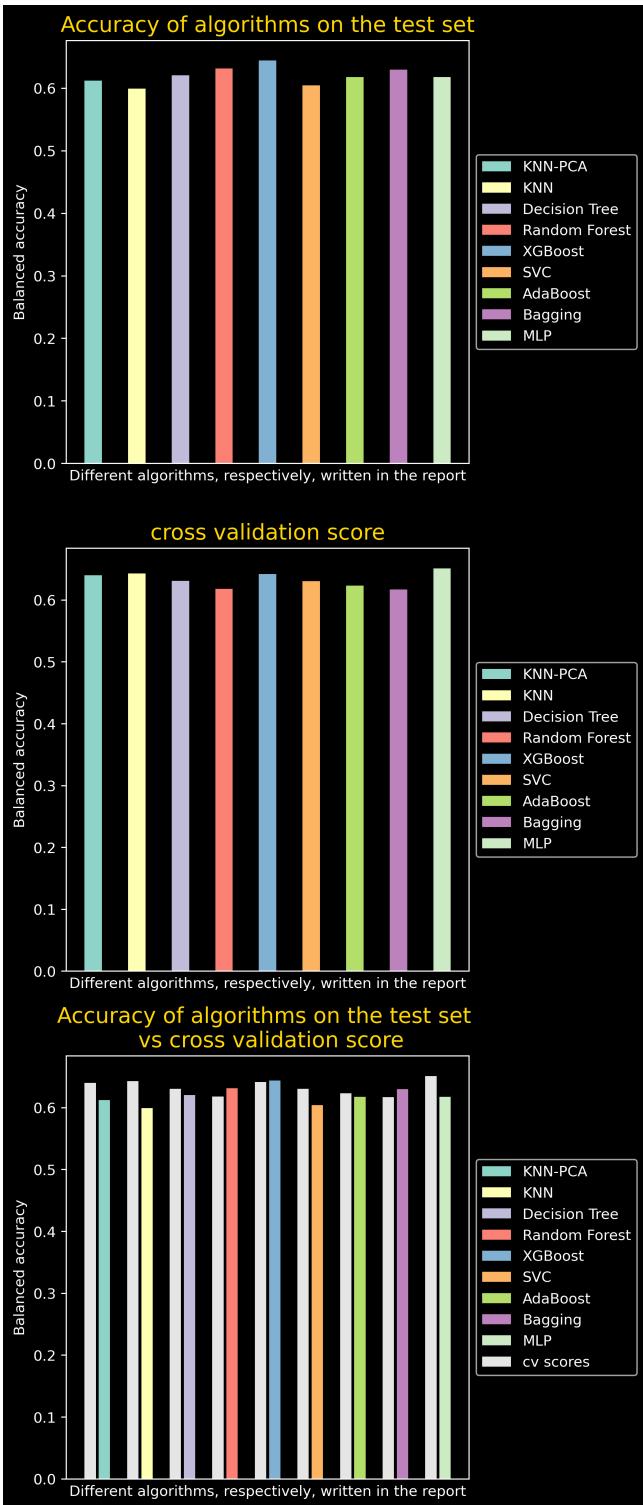


Figure 12: Diagrams of the 5th dimension of the labels vectors[4]

related to KNN-PCA and Bagging. For this reason, KNN-PCA and KNN have performed well or well, almost wherever Bagging has performed poorly. Interestingly, after KNN-PCA and KNN, the highest correlation with a high rate of 96% is related to Random forest and Bagging algorithms. So with a bit of carelessness, I can claim everything that I said about Bagging, about the Random forest algorithm.[Refer to Figure 16]

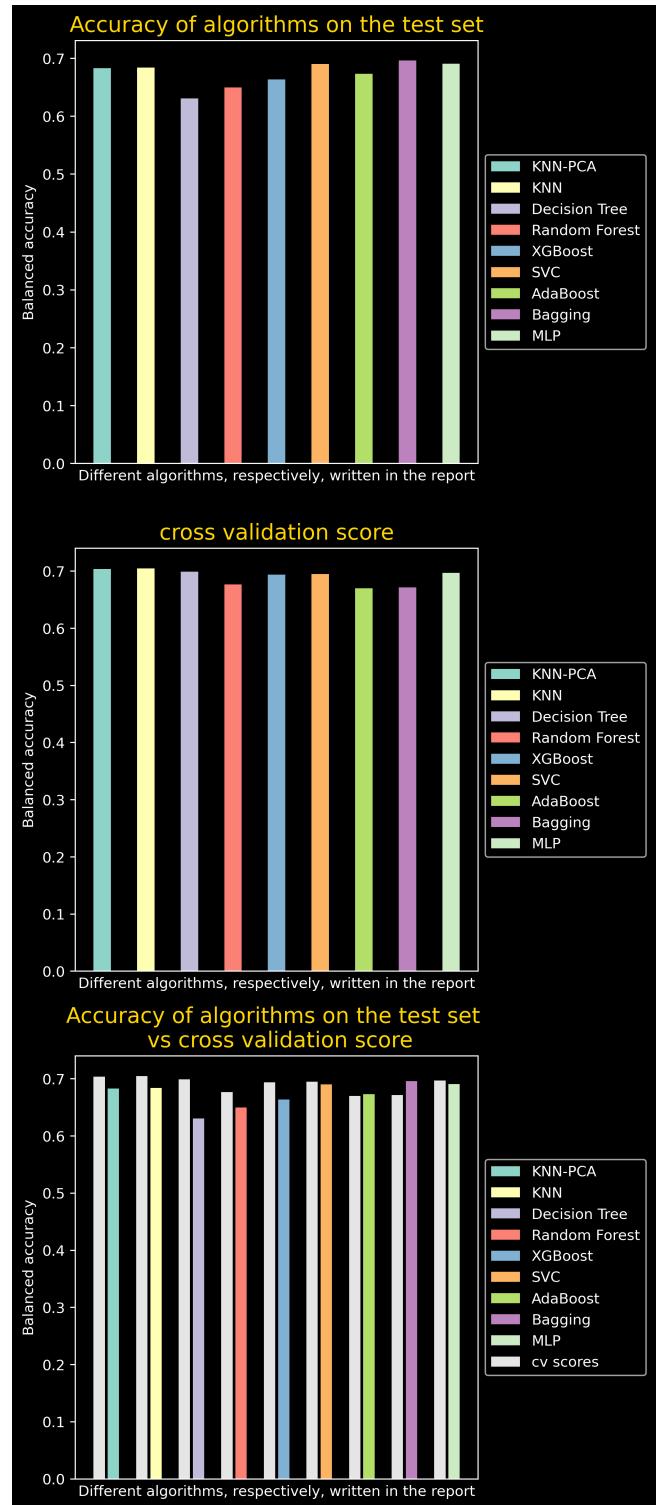


Figure 13: Diagrams of the 6th dimension of the labels vectors[4]

## 7.9 9th dimension of the labels vector

The results of the algorithms on the label vector's ninth dimension, which represents the toxicity of SR-ATAD5, can be seen in Figure 17. As shown in Figure 22, this dimension has one of the highest variances of algorithm accuracy compared to other dimensions. In this dimension, KNN-PCA has obtained the highest score based on the balanced accuracy measurement system. As expected, Bagging was the least accurate.

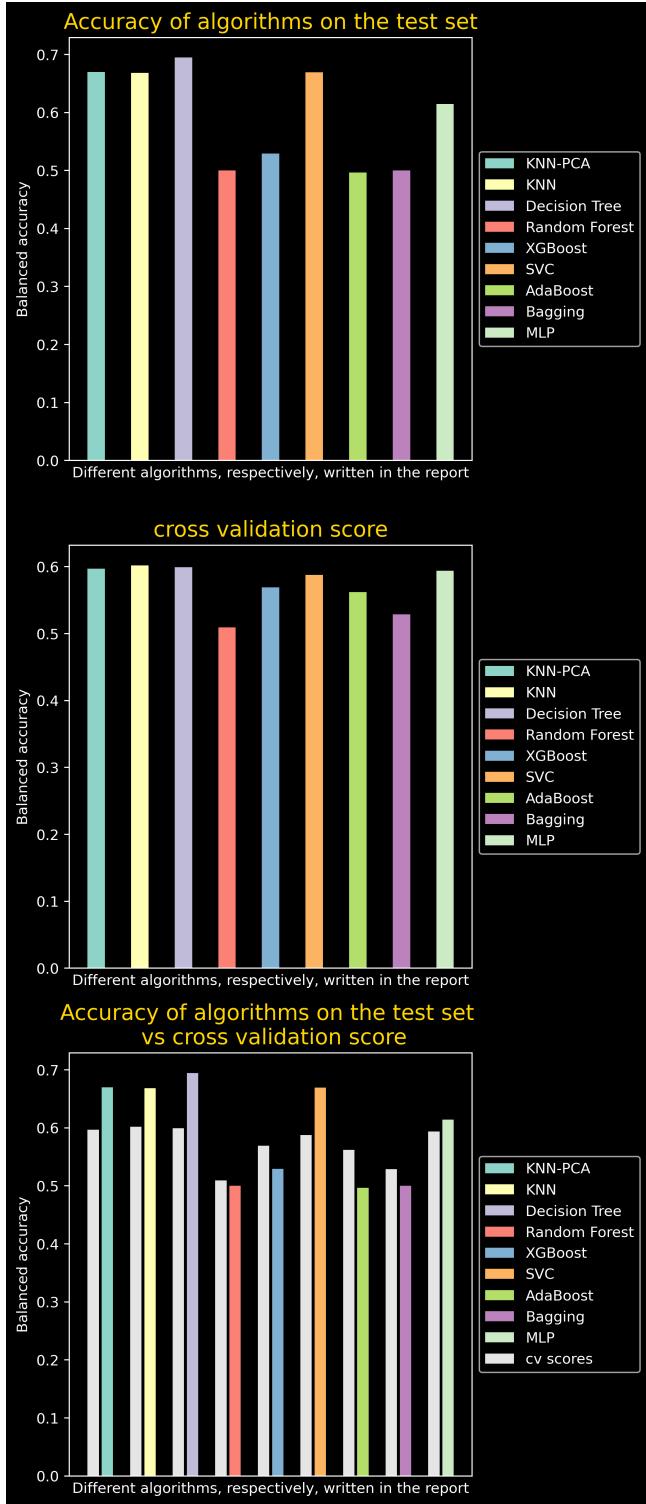


Figure 14: Diagrams of the 7th dimension of the labels vectors[4]

But given the high correlation between the labels predicted by KNN-PCA and KNN, the question is, did KNN also perform well? As shown in Figure 10, the ninth dimension shows the greatest advantage of KNN-PCA over KNN. Figure 17 also shows a small amount of overfitting in the KNN. Of course, the hyperparameters are selected in exactly the same in these algorithms, but KNN-PCA correctly predicts about 25% more True labels than KNN.

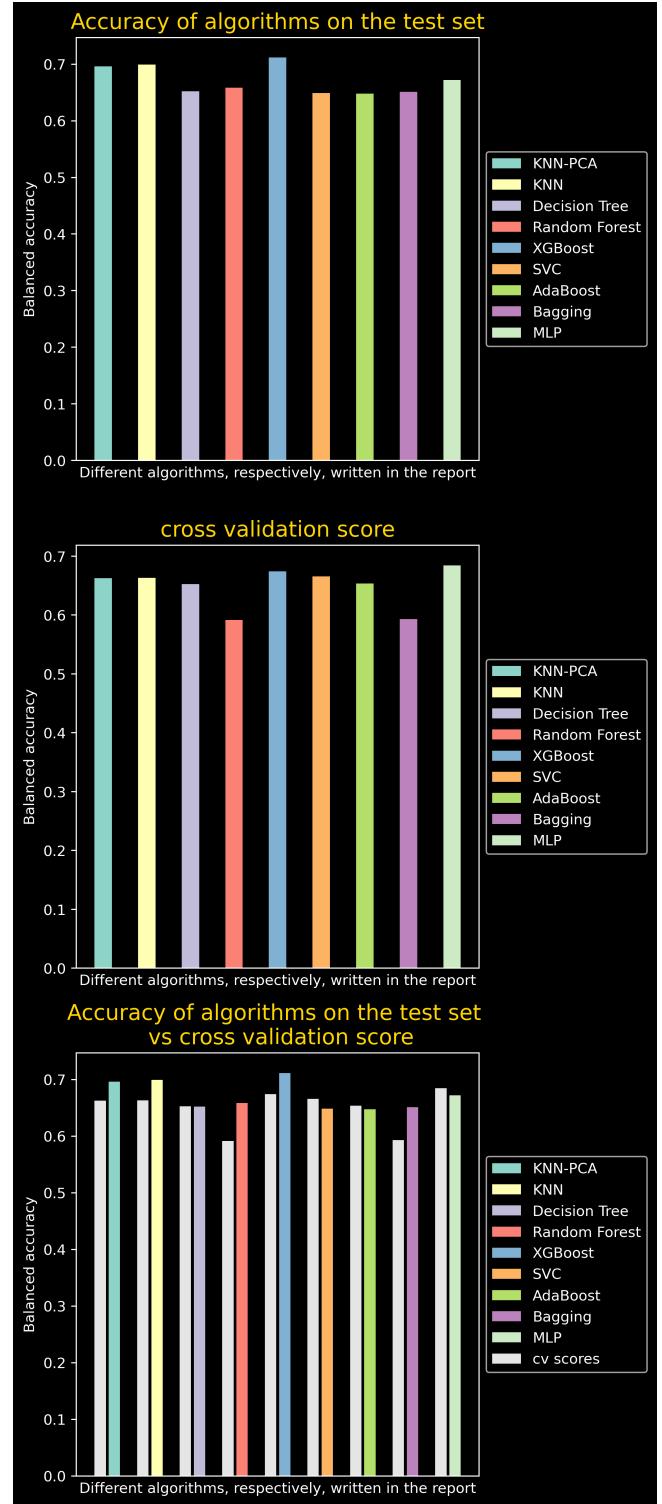


Figure 15: Diagrams of the 8th dimension of the labels vectors[4]

## 7.10 10th dimension of the labels vector

The results of performing algorithms on the tenth dimension are very similar to the ninth dimension. In this dimension, the label vectors representing SR-HSE toxicity can be seen in Figure 18. As shown in Figure 22, this dimension, like the ninth dimension, has a significant variance in the accuracy of the algorithms. KNN-PCA has obtained the highest score based on the balanced accuracy measurement system in the

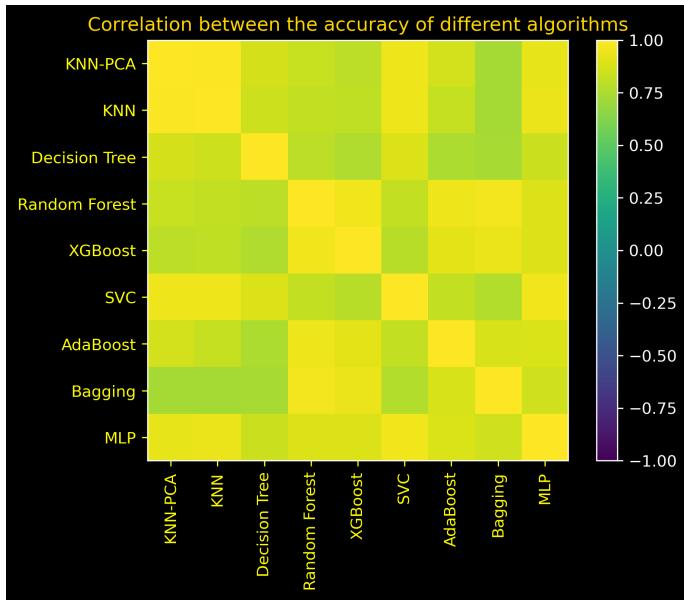


Figure 16: Correlation between the accuracy of different algorithms[4]

tenth dimension. As expected, the result of Bagging was not desirable, but the least accurate is related to the Random forest algorithm.

In Figure 23, you can see the accuracy range of the algorithm in these 12 dimensions (average, minimum and maximum). As the figure shows, the Random forest algorithm has a relatively high variance in balanced accuracy on these 12 dimensions. This algorithm had the highest score in predicting NR-AR toxicity, while it had the lowest accuracy here.

Also, despite the fact that the KNN-PCA hyperparameter<sup>11</sup> has been selected in this dimension (as well as in many other dimensions) in such a way that it seems to lead to overfitting, our test results show that we are not trapped by overfitting at all.

### 7.11 11th dimension of the labels vector

In the eleventh dimension of the feature vector, we refer to the toxicity of SR-MMP [Refer to Figure 19]. The accuracy of the algorithms has been almost the same. The highest accuracy on the test set was related to the XGBoost algorithm. Unfortunately, the test set result was slightly lower than the cv score. But this difference is less than what we worry about overfitting. While the highest cv score was for MLP, the final accuracy on the test set was significantly different from the cv score, which makes me motivated to choose XGBoost.

The lowest accuracy of both test set results and the cv score result are related to the bagging algorithm. After Bagging, the weakest performance was associated with the Random forest algorithm. Of course, I emphasize again that the maximum difference in data accuracy was 4%, which means that all algorithms have performed approximately the same. That's why Random forest may be a good option because it works with fewer features, which makes it easier for us to collect data. But in terms of accuracy, XGBoost was a little better than the rest.

### 7.12 12th dimension of the labels vector

Figures 20 and 22 show us that not all algorithms in the twelfth dimension are alike. However, mathematically the highest balanced accuracy belongs to KNN. I want to break with tradition and choose MLP as the best model here. The MLP algorithm has about 0.2% less balanced accuracy than KNN, which is not a significant number at all. Because the test score of the MLP algorithm is much better than the cv score, I can firmly claim that in MLP, we taught the model even better than KNN.

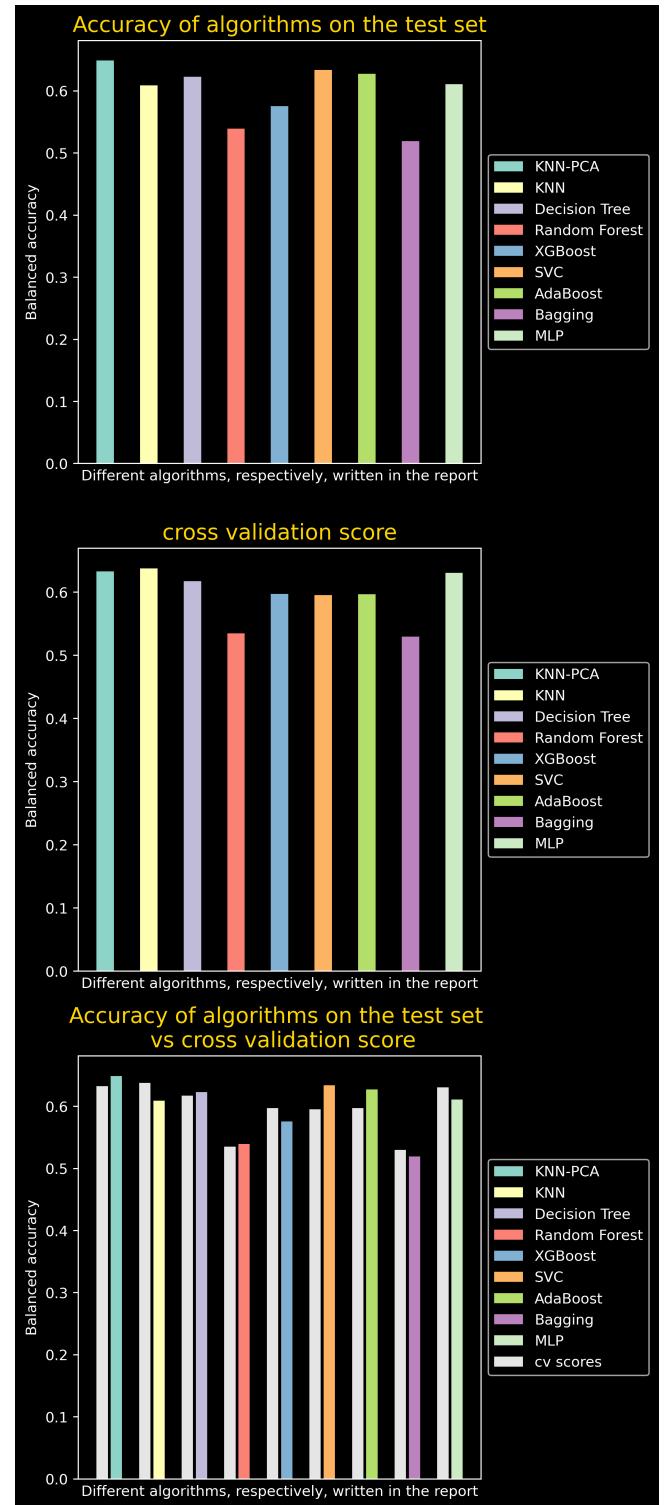


Figure 17: Diagrams of the 9th dimension of the labels vectors[4]

In SR-p53-related labels, such as the SR-MMP dimension, the weakest performance is in Bagging and Random forest, respectively. But in this dimension, unlike the eleventh dimension, the difference between these two algorithms and the best algorithms is much more than what can be ignored.

In general, I think it is easier for us to choose when the algorithms differ from the results of different algorithms. Because when the

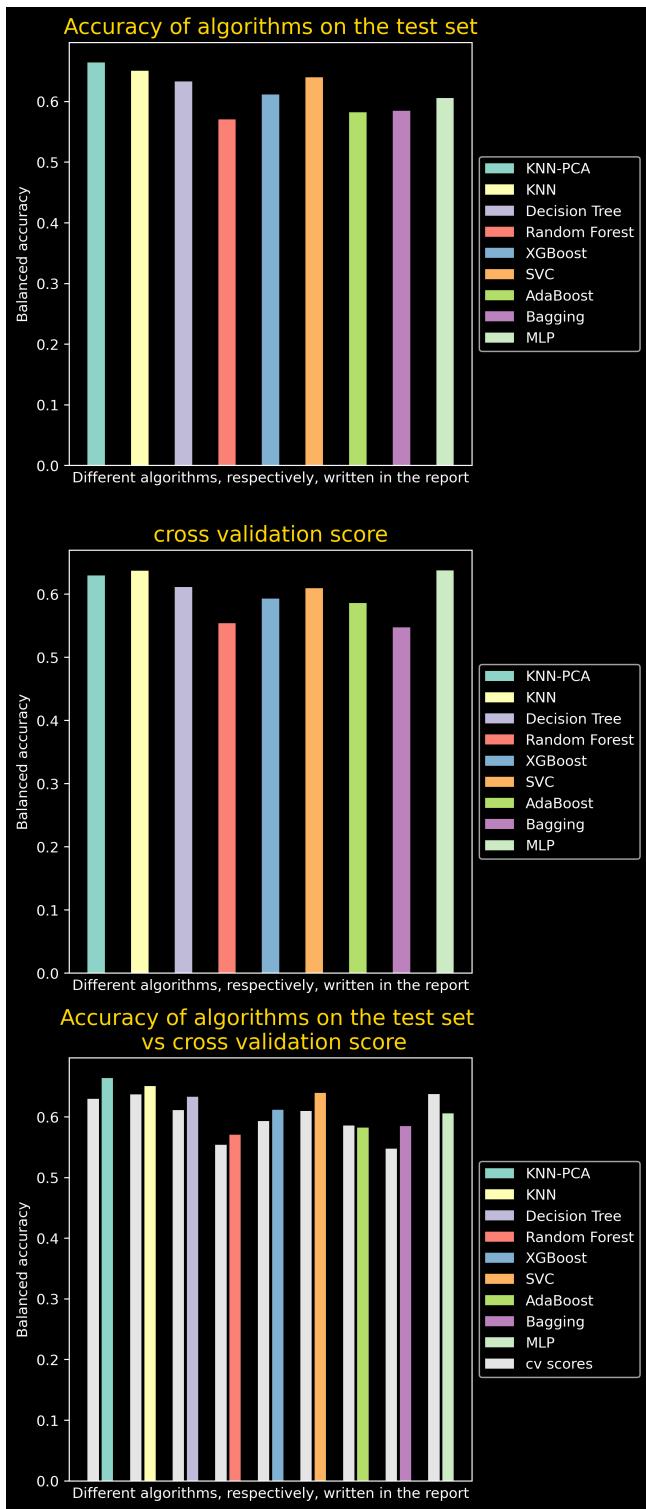


Figure 18: Diagrams of the 10th dimension of the labels vectors[4]

difference between the results of different algorithms is slight, we have to consider various parameters. For example, one algorithm maybe 0.002 percent weaker than another if it has hundreds of times simpler models or takes much less time to run this algorithm. Therefore, when one algorithm performs better than other algorithms with a big difference, we are less involved in these analyzes. As a result, the probability of our analytical error also decreases significantly.

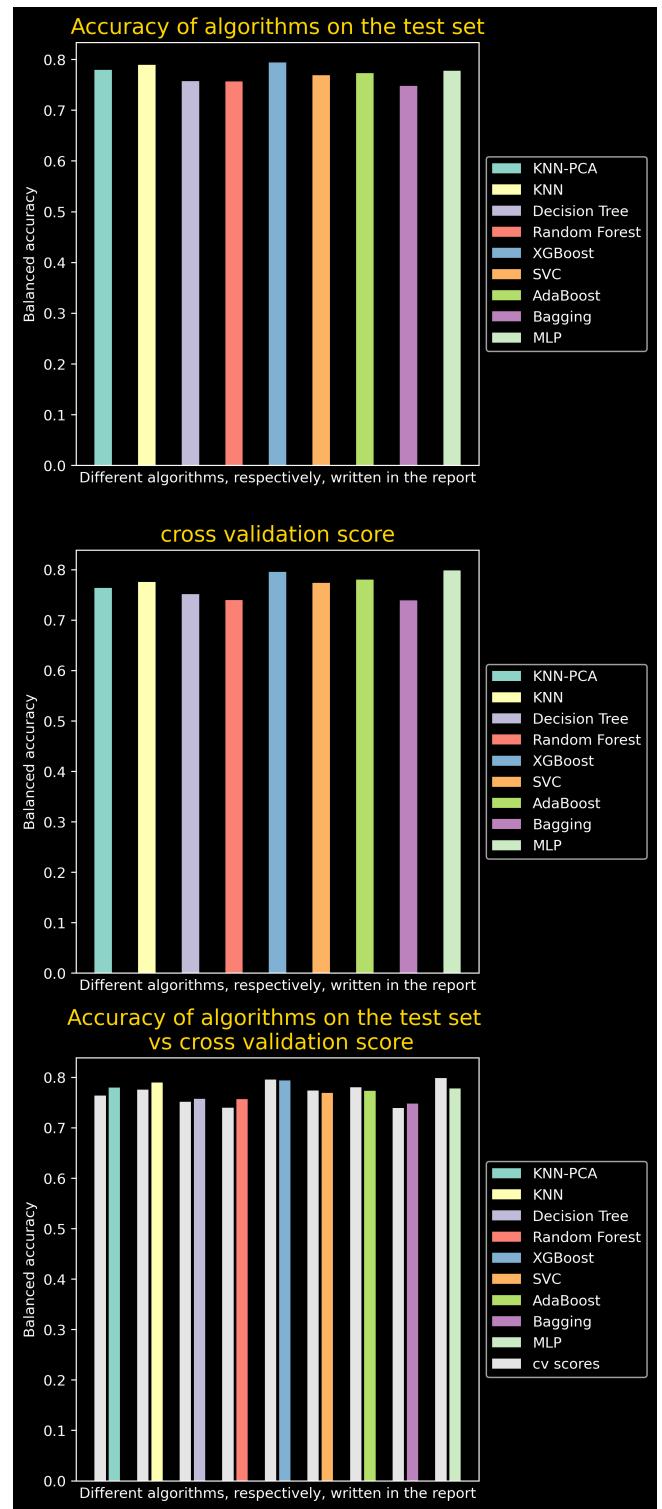


Figure 19: Diagrams of the 11th dimension of the labels vectors[4]

## 8 CONCLUSION

### 8.1 Initial conclusion

First, I want to briefly rewrite the selected algorithms in each dimension and the accuracy of the information associated with them. Then I will provide an overview of which algorithms we used the most and which algorithms we did not use. Also, as a multi-label project, I must specify the final accuracy.

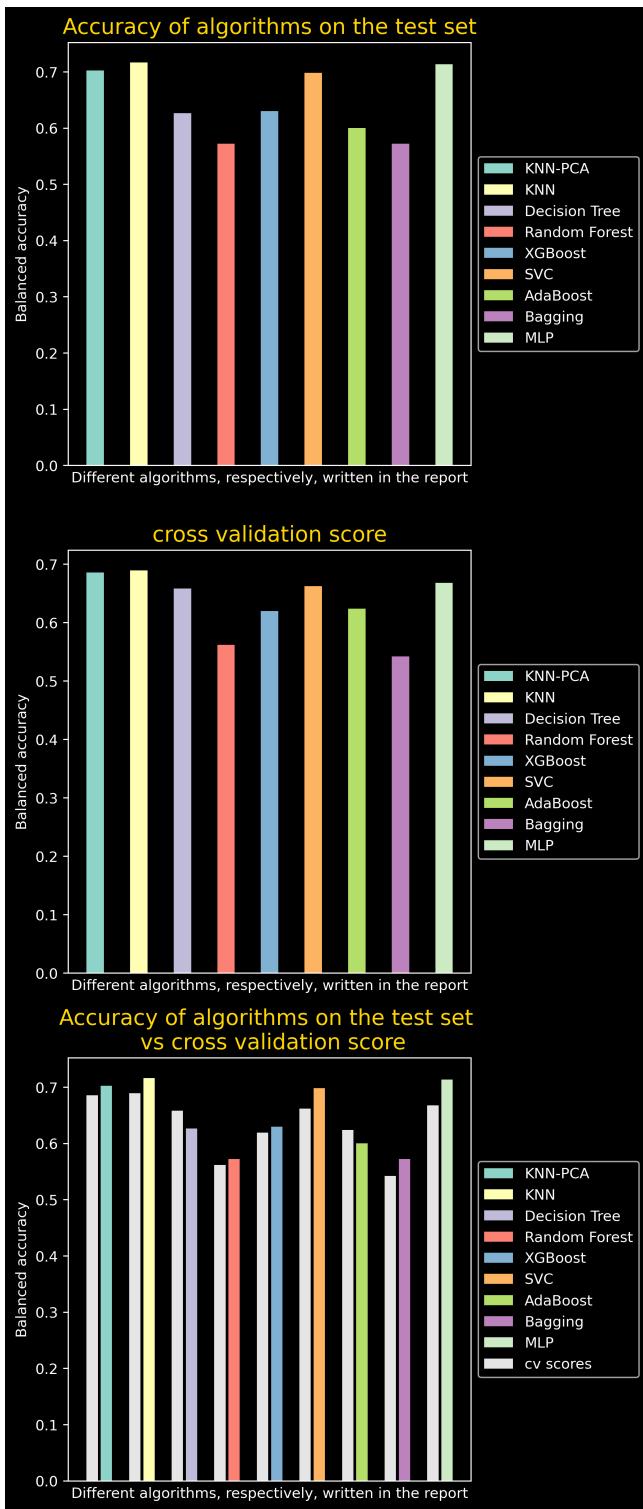


Figure 20: Diagrams of the 12th dimension of the labels vectors[4]

**2. NR-AR-LBD**

- Selected algorithm: Random forest
- Hyperparameters: criterion: gini, max\_depth: 15, max\_features: auto, n\_estimators: 700
- Balanced accuracy: 75.66067901815103%

**3. NR-AhR**

- Selected algorithm: KNN
- Hyperparameters: metric: manhattan, n\_neighbors: 1
- Balanced accuracy: 76.94407169013123%

**4. NR-Aromatase**

- Selected algorithm: KNN-PCA
- Hyperparameters: metric: cosine, n\_neighbors: 1
- Balanced accuracy: 76.7794028031688%

**5. NR-ER**

- Selected algorithm: XGBoost
- Hyperparameters: eval\_metric: mlogloss, max\_depth: 3, n\_estimators: 2000
- Balanced accuracy: 64.41829042810411%

**6. NR-ER-LBD**

- Selected algorithm: Bagging
- Hyperparameters: max\_features: 0.8, max\_samples: 0.5, n\_estimators: 100
- Balanced accuracy: 69.62051709758133%

**7. NR-PPAR-gamma**

- Selected algorithm: Decision Tree
- Hyperparameters: criterion: gini, max\_depth: 28
- Balanced accuracy: 69.46102112342521%

**8. SR-ARE**

- Selected algorithm: XGBoost
- Hyperparameters: eval\_metric: mlogloss, max\_depth: 5, n\_estimators: 2000
- Balanced accuracy: 71.15723782390448%

**9. SR-ATAD5**

- Selected algorithm: KNN-PCA
- Hyperparameters: metric: cosine, n\_neighbors: 1
- Balanced accuracy: 64.88888888888889%

**10. SR-HSE**

- Selected algorithm: KNN-PCA
- Hyperparameters: metric: manhattan, n\_neighbors: 1
- Balanced accuracy: 66.41911069063386%

**11. SR-MMP**

- Selected algorithm: XGBoost

**1. NR-AR**

- Selected algorithm: Random forest
- Hyperparameters: criterion: gini, max\_depth: 15, max\_features: auto, n\_estimators: 700
- Balanced accuracy: 75.66067901815103%

- Hyperparameters: eval\_metric: mlogloss, max\_depth: 3, n\_estimators: 700
- Balanced accuracy: 79.39393939393939%

## 12. SR-p53

- Selected algorithm: MLP<sup>12</sup>
- Hyperparameters: activation: tanh , hidden\_layer\_sizes: 500
- Balanced accuracy: 71.39478072046221%

As you can see in the histogram of Figure 21, the frequency of selection of KNN-PCA and XGBoost algorithms was three times. The Random Forest algorithm was then selected twice to predict the two dimensions of the label vector. The KNN, Bagging, Decision Tree, and MLP algorithms were in third place, each selected once. Also, the SVC and AdaBoost algorithms were never selected.

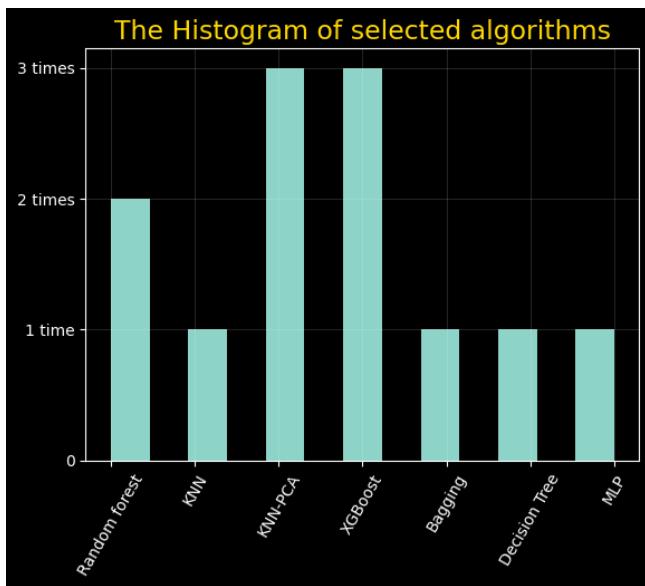


Figure 21: The Histogram of selected algorithms[4]

I also included Figure 22 so that you can see the performance of each algorithm in different dimensions of the labels vector in one frame.

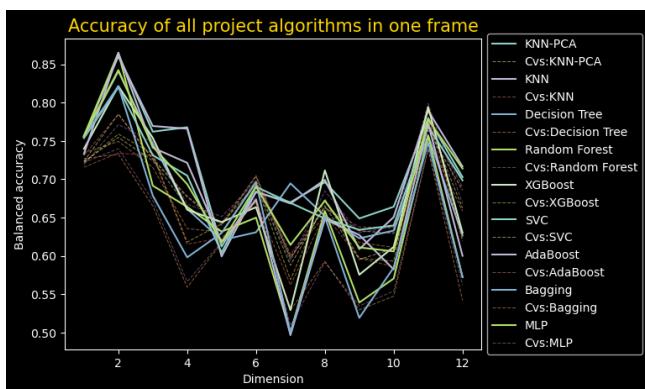


Figure 22: Accuracy of all project algorithms in one frame[4]

Since Figure 22 might seem a bit crowded, I prepared Figure 23 to evaluate the algorithms better so that you can see the accuracy range of each algorithm and the average accuracy of each algorithm separately.

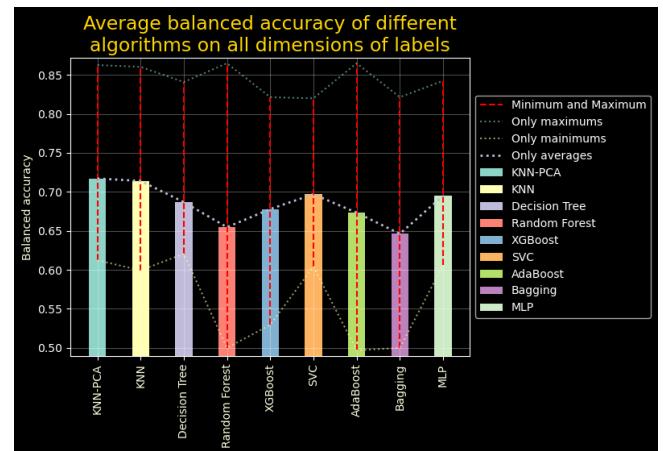


Figure 23: Average balanced accuracy of different algorithms on all dimensions of labels[4]

I want to say something informal in this part. Of course, the whole project had to be multitasking, and all 12 dimensions had to be applied to the data simultaneously without assuming independence. But according to the professor's orders, I decided to implement the project as above. Of course, we examined the correlation between labels, and the assumption of label independence is not very wrong.

However, if we were to study 12 dimensions simultaneously, we would probably have much better options of balanced accuracy to choose as a meter. Although, I want to introduce you to a very informal criterion based on balanced accuracy for evaluating the whole project, which was just made up by my own mind. After thinking about it, I saw many weaknesses, but I preferred to write this criterion in my report.

My goal from the beginning was to predict a 12-dimensional vector. I have said before that I use balanced accuracy to evaluate my algorithms. But I would like to have a number as the final accuracy of these algorithms on the whole project. First, this number must be good. This means that it cannot be interpreted as a liar. Secondly, it should be interpretable.

Since we had a confusion matrix for each selected algorithm, I guessed that summing these confusion matrices might give us good information. It should not be left unmentioned that average balancing is not a suitable parameter because the distribution of zeros and ones and the number of test set members in each label are different. I finally arrive at an interpretable matrix by summing the confusion matrices. This matrix explicitly answers the following questions.

1. How many True labels did I guess correctly?
2. How many False labels did I guess correctly?
3. How many True labels did I guess wrong?
4. How many False labels did I guess wrong?

Finally, I can calculate the balanced accuracy of the matrix and reach an excellent meter to measure the overall accuracy of the project.

$$\text{Confusion matrix} = \begin{bmatrix} 6880 & 175 \\ 293 & 279 \end{bmatrix}$$

$$\text{Balanced accuracy} = 73.14785674991204\%$$

In other words, I correctly predicted 97.51948972360028% of zeros in the total whole of the project. I also predicted 48.776223776223776% of the ones correctly. Clearly, our main challenge throughout the project has been to anticipate the ones.

There may be many reasons why this criterion is not good. For example, to say that we need a measure that examines every twelve-dimensional vector as a package of 12 binaries. Suppose we have a test set of 500 members. Our labels will be in the form of the following matrix.

$$\text{Labels : } \begin{bmatrix} 0 & \dots & 1 \\ 1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix}_{500 \times 12}, \text{ Prediction : } \begin{bmatrix} 0 & \dots & 1 \\ 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1 \end{bmatrix}_{500 \times 12}$$

My goal is to know how similar these two matrices are. As a result, we can define a confusion matrix for this matrix and know how accurate the prediction matrix was. Summing the confusion matrices I mentioned earlier is exactly the same as the confusion matrix obtained by comparing the above two matrices. I got the basic idea from studying a response in stackexchange.[19]

## 8.2 Final project conclusion

While editing the project report, I noticed that I could not select the final model using the test set. In competitions, to be fair, the test dataset is not revealed until you submit the final trained model.

Therefore, I decided to select the best models based on the cv score at the end of the report and compare the newly chosen models with the previous ones selected based on the test set results.

### 1. NR-AR

- Selected algorithm: Decision Tree
- CV score<sub>based on the balanced accuracy</sub>: 73.43343187519518%
- Balanced accuracy on the test set: 75.51490642339886%

### 2. NR-AR-LBD

- Selected algorithm: KNN-PCA
- CV score<sub>based on the balanced accuracy</sub>: 78.5002324569213%
- Balanced accuracy on the test set: 86.25992731188586%

### 3. NR-AhR

- Selected algorithm: MLP
- CV score<sub>based on the balanced accuracy</sub>: 74.85781227406562%
- Balanced accuracy on the test set: 74.25910341847295%

### 4. NR-Aromatase

- Selected algorithm: KNN
- CV score<sub>based on the balanced accuracy</sub>: 67.81417541315211%
- Balanced accuracy on the test set: 76.59658744667885%

### 5. NR-ER

- Selected algorithm: MLP
- CV score<sub>based on the balanced accuracy</sub>: 65.13103184676162%
- Balanced accuracy on the test set: 61.77049258047693%

### 6. NR-ER-LBD

- Selected algorithm: KNN
- CV score<sub>based on the balanced accuracy</sub>: 70.49780678881433%
- Balanced accuracy on the test set: 68.39727550736725%

### 7. NR-PPAR-gamma

- Selected algorithm: KNN
- CV score<sub>based on the balanced accuracy</sub>: 60.20852568911459%
- Balanced accuracy on the test set: 66.84190584446338%

### 8. SR-ARE

- Selected algorithm: MLP
- CV score<sub>based on the balanced accuracy</sub>: 68.45563149901453%
- Balanced accuracy on the test set: 67.20820331931444%

### 9. SR-ATAD5

- Selected algorithm: KNN
- CV score<sub>based on the balanced accuracy</sub>: 63.77311937912506%
- Balanced accuracy on the test set: 60.88888888888888%

### 10. SR-HSE

- Selected algorithm: MLP
- CV score<sub>based on the balanced accuracy</sub>: 63.78664632826452%
- Balanced accuracy on the test set: 60.60075685903501%

### 11. SR-MMP

- Selected algorithm: MLP
- CV score<sub>based on the balanced accuracy</sub>: 79.8803902979208%
- Balanced accuracy on the test set: 77.80532598714417%

### 12. SR-p53

- Selected algorithm: KNN
- CV score<sub>based on the balanced accuracy</sub>: 68.9557702999988%
- Balanced accuracy on the test set: 71.66039784404204%

With these interpretations, you can see the frequency of selected algorithms based on the Cv score in Figure 24.

I want to take this opportunity to try the metrics I suggested to evaluate the whole project here. We expect the accuracy to be a bit lower here, but we should not take it too seriously from the previous confusion matrix. I calculated the cumulative confusion matrix for the selected models based on the cross-validation score.

$$\text{Confusion matrix} = \begin{bmatrix} 6704 & 210 \\ 322 & 257 \end{bmatrix}$$

$$\text{Balanced accuracy} = 70.67477916449965\%$$

In other words, I correctly predicted 96.96268440844663% of zeros in the total whole of the project. I also predicted 44.38687392055268% of the ones correctly.

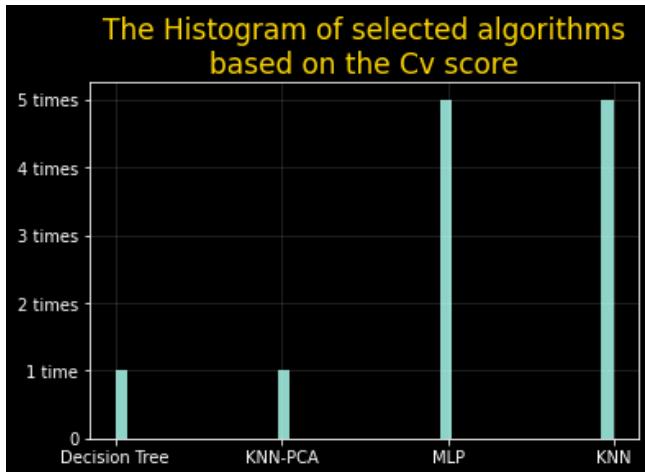


Figure 24: The Histogram of selected algorithms based on the Cv score[4]

As you can see, we did not use the test set to re-select the algorithm. So it was apparent that our criterion called cumulative confusion matrix should give a lower score than the whole project, which it did. This shows that this indicator does not seem very illogical.

The entire project run time (excluding research and programming) for me lasted about four consecutive days on a 48-core computer (each core was 2.4 GHz).

## NOTES

<sup>1</sup>initial.dataset.dropna().describe().iloc[0,0] command in the Python environment after reading the label file can return us the number.

<sup>2</sup>**Principal component analysis (PCA)** is a technique for reducing the dimensionality of such datasets, increasing interpretability but at the same time minimizing information loss. It does so by creating new uncorrelated variables that successively maximize variance.[20]

<sup>3</sup>According to the group's agreement, the random number 1234 has been selected for train\_test\_split

<sup>4</sup>**GridSearchCV** implements a **fit** and a **score** method. It also implements **score\_samples**, **predict**, **predict\_proba**, **decision\_function**, **transform** and **inverse\_transform** if they are implemented in the estimator used.[5]

<sup>5</sup>TP:True Positive, P:Positive, TN:True Negative, N:Negative

<sup>6</sup>FP:False Positive ,FN: False Negative

<sup>7</sup>for example, when the training set is the entire Web and distances are based on links

<sup>8</sup>Independent and identically distributed

<sup>9</sup>In statistics, overfitting is "the production of an analysis that corresponds too closely or exactly to a particular set of data, and may therefore fail to fit additional data or predict future observations reliably". An overfitted model is a statistical model that contains more parameters than can be justified by the data. The essence of overfitting is to have unknowingly extracted some of the residual variation (i.e., the noise) as if that variation represented underlying model structure.

<sup>10</sup>As you know, the best hyper-parameters are selected based on this parameter.

<sup>11</sup>In machine learning, a hyperparameter is a parameter whose value is used to control the learning process. By contrast, the values of other parameters (typically node weights) are derived via training.

<sup>12</sup>Multi Layer Perceptron

## ACKNOWLEDGEMENTS

I would like to thank my applied machine learning professor (Dr. Hajiabolhasan) and all teacher assistants who gave me a golden opportunity to work on this project. I wish these loved ones good health and success.

## REFERENCES

- [1] Shai Shalev-Shwartz and Shai Ben-David. Understanding Machine Learning - From Theory to Algorithms. Cambridge University Press, 2014.
- [2] Li Di and Edward H. Kerns. Chapter 17 - toxicity. In Li Di and Edward H. Kerns, editors, Drug-Like Properties (Second Edition), pages 251–262. Academic Press, Boston, second edition edition, 2016.
- [3] hulab.rxnfinder.org. Convert smiles to image, 2019.
- [4] J. D. Hunter. Matplotlib: A 2d graphics environment. Computing in Science & Engineering, 9(3):90–95, 2007.
- [5] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12:2825–2830, 2011.
- [6] Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming with NumPy. Nature, 585(7825):357–362, September 2020.
- [7] Wes McKinney. Data Structures for Statistical Computing in Python. In Stéfan van der Walt and Jarrod Millman, editors, Proceedings of the 9th Python in Science Conference, pages 56 – 61, 2010.
- [8] Jr Mohammed J. Zaki, Wagner Meira. Data Mining and Machine Learning-Fundamental Concepts and Algorithms. Cambridge University Press, 2020.
- [9] Joydwip Mohajon. Confusion matrix for your multi-class machine learning model, 2020.
- [10] Xindong Wu, Vipin Kumar, J. Ross Quinlan, Joydeep Ghosh, Qiang Yang, Hiroshi Motoda, Geoffrey J. McLachlan, Angus F. M. Ng, B. Liu, Philip S. Yu, Zhi-Hua Zhou, Michael S. Steinbach, David J. Hand, and Dan Steinberg. Top 10 algorithms in data mining. Knowledge and Information Systems, 14:1–37, 2007.
- [11] Tianqi Chen and Carlos Guestrin. XGBoost: A scalable tree boosting system. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’16, pages 785–794, New York, NY, USA, 2016. ACM.
- [12] Corinna Cortes and Vladimir Naumovich Vapnik. Support-vector networks. Machine Learning, 20:273–297, 2004.
- [13] Wikipedia contributors. Adaboost Wikipedia, the free encyclopedia, 2022. [Online; accessed 13-Feb-2022].
- [14] Wikipedia contributors. Bootstrap Wikipedia, the free encyclopedia, 2021. [Online; accessed 13-Feb-2022].
- [15] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. The Elements of Statistical Learning. Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.
- [16] Frank Rosenblatt. Principles of neurodynamics. perceptrons and the theory of brain mechanisms. American Journal of Psychology, 76:705, 1963.
- [17] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning internal representations by error propagation. In David E. Rumelhart and James L. McClelland, editors, Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Volume 1: Foundations, pages 318–362. MIT Press, Cambridge, MA, 1986.
- [18] G. Cybenko. Approximation by superpositions of a sigmoidal function. Mathematics of Control, Signals, and Systems (MCSS), 2(4):303–314, December 1989.
- [19] rapaio. How to get an aggregate confusion matrix from n different classifications, 2014. [Online; accessed 10-Jun-2014].
- [20] Jorge Cadima Ian T Jolliffe. Principal component analysis: a review and recent developments. Philosophical transactions. Series A, Mathematical, physical, and engineering sciences, (7825):374, April 2016.