

Hiv 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, Logistic regression, Boosting, Bagging, AdaBoost, Ensemble, Random forest, Decision tree, Python, Balanced accuracy, Confusion matrix

ABSTRACT

The following project is a machine learning project written in Python. This project's dataset describes the formula of chemical molecules of 41,127 different drugs. I should mention that the chemical properties of each of them were available to us with a 200-dimensional normalized vector. Using understanding machine learning from theory to algorithms [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

Contents	1
1 Introduction	1
2 Description of raw data	2
2.1 Labels distribution	2
2.2 Missing data	2
2.3 Algebraic dimension of features	2
3 Data separation	2
4 dimensionality reduction	2
5 Algorithms	3
5.1 K Nearest Neighbors	3
5.2 Decision Tree	3
5.3 Random Forest	4
5.4 XGBoost	4
5.5 C-Support Vector Classification	5
5.6 Logistic Regression	5
5.7 AdaBoost	5
5.8 Bagging	6
6 Evaluation	6
References	8

PRE-INTRODUCTION

I'm so sorrowful that I have to start my report with shocking statistics. My motivation for choosing this project was that this disease has led to the death of many victims. This project is definitely too small to have the slightest impact. But I think the essential message of this project is the

statistics ahead, which motivated me to work on this project. I hope the definitive cure will be discovered as soon as possible.

- 1.5 million people became newly infected with HIV in 2020.
- 680,000 people died from AIDS-related illnesses in 2020.
- 79.3 million people have become infected with HIV since the start of the epidemic.¹
- 36.3 million people have died from AIDS-related illnesses since the start of the epidemic.²
- Studies from England and South Africa have found that the risk of dying from COVID-19 among people with HIV was double that of the general population. [2]

1 INTRODUCTION

My data contains two files named `hiv_global_cdf_rdkit` and `hiv`. 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]

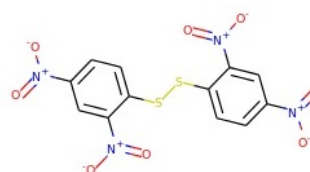


Figure 1: The chemical formula of the eighth data

The label file is relatively straightforward and unambiguous, but the feature file contains a significant amount of empty information. Of the 41,127 data, 39,629 are wholly unambiguous. I should note that data labels are highly imbalanced. I will discuss each of these points in detail in the following parts. But I must now point out that this data has two main challenges.

1. The relatively high number of data will lead to a large volume of calculations
2. Data imbalance may be an area of expertise that I am certainly not specialized in.

2 DESCRIPTION OF RAW DATA

2.1 Labels distribution

As I mentioned before, we are dealing with highly imbalanced data. For further explanation, I must say that out of 41,127 data; only 1443 data are labeled 1. This means that only 3.508% of the data are labeled 1, and 96.491% are labeled zero. For a better understanding, Figure 2 shows this imbalanced distribution.

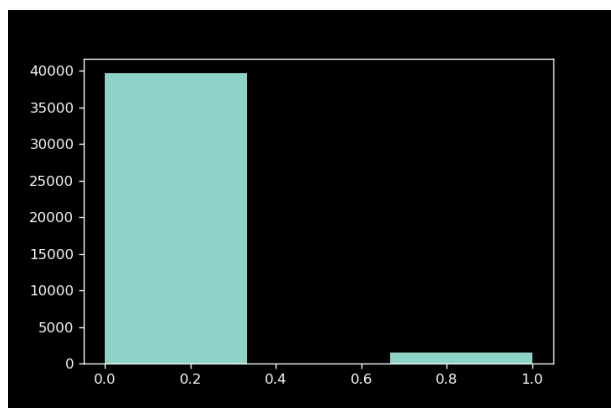


Figure 2: Labels distribution

2.2 Missing data

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.[4]

But a significant challenge is missing data. Is this data damaged as same as the same original distribution? No! Out of 1498 missing data, 172 of data have label 1, which is approximately 11.5% of the total missing data. In other words, if we want to discard all the damaged data, the balance of the remaining information is not the same as reality. If the damaged data is discarded, only 3.207% of the remaining data will be labeled as 1. This means that more than 11.9% of the tagged one information is somehow damaged. To solve this problem, I separated the data with label one then replaced the missing feature with the average of that column.³ 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 #Importing essential packages
5 initial_Features=pd.read_csv('hiv_global_cdf_rdkit.csv')
6 initial_dataset=pd.read_csv('hiv.csv')
7 #Reading data
8 ones_array=[]
9 only_ones_inf=initial_Features
10 for i in range(0,41127):
11     if initial_dataset['HIV_active'][i]>0:
12         ones_array.append(i)
13     else:
14         only_ones_inf=only_ones_inf.drop(i)
15 #Separate data labeled 1
16 only_ones_inf_filled_mean=only_ones_inf.fillna(
17     only_ones_inf.mean())
18 #Fill in the blank data with the average of the same
19 #property (only for data that has a label.)
20 only_zeros_inf=initial_Features
21 for i in ones_array:
```

```
20     only_zeros_inf=only_zeros_inf.drop(i)
21 #Separate data labeled 0
22 secondary_Features=pd.concat([
23     only_ones_inf_filled_mean,only_zeros_inf]).
24     sort_index().dropna()
25 #Define new data (without deleting data that is
26 #labeled one)
27 secondary_dataset=pd.read_csv('hiv.csv')
28 tmp=pd.concat([only_ones_inf_filled_mean,
29     only_zeros_inf]).sort_index().isna().any(axis=1)
30 np.size(tmp)
31 for i in np.arange(0,np.size(tmp)):
32     if tmp[i]==True:
33         secondary_dataset=secondary_dataset.drop(i,
34             axis=0)
35 #Syncing secondary labels according to secondary data
```

Listing 1: Deal with missing data [5] [6]

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 192 algebraic dimensions.[Refer to Listing 2] 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.

```
1 secondary_Features=secondary_Features.loc[:,
2     secondary_Features.apply(pd.Series.nunique) != 1]
```

Listing 2: Remove ineffective features

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 PCA⁴ 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.⁵ I also intended to allocate 10% of the data volume to the test set.

```
1 from sklearn.model_selection import train_test_split
2 X_training_data, X_test, y_training_data, y_test=
3     train_test_split(secondary_Features,
4         secondary_dataset, stratify=secondary_dataset['
5         HIV_active'],test_size=0.10,random_state=1234)
```

Listing 3: using train test split [4]

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 [7]

I also looked at the data and found that PCA on my data, with 77 principal components, can maintain 95.0935% of the data variance. So, I made this dimensional reduction for some algorithms, as shown in Listing 4

```
1 from sklearn.decomposition import PCA
2 pca = PCA(n_components=77)
3 principalComponents = pca.fit_transform(
4     X_training_data)
5 X_training_data_PCA = pd.DataFrame(data =
6     principalComponents)
```

Listing 4: PCA [4]

5 ALGORITHMS

In this section, I will run various algorithms on the data. I will give a specific explanation about each algorithm. Then you will see the result of running the algorithm. I have to say that I ran the algorithms using gridsearchcv and chose the best hyperparameters.⁶

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.

$$\text{Balanced Accuracy} = \frac{1}{2} \left(\frac{TP}{P} + \frac{TN}{N} \right) \quad 7$$

Also, in the algorithms results, you will encounter a matrix called confusion matrix, which is defined as in Figure 3.

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

Figure 3: confusion matrix[8]⁸

5.1 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 fast⁹. [1] You can see the execution of KNN on 77-dimensional data (after PCA application) in Listing 5.

```
1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.model_selection import GridSearchCV
3 from sklearn.metrics import confusion_matrix,
4     plot_confusion_matrix
5 knn = KNeighborsClassifier()
6 param_grid={'n_neighbors':np.arange(1,100,2),'metric':
7     :['euclidean','cosine','manhattan']}
8 grid=GridSearchCV(knn,param_grid=param_grid,cv=10,
9     scoring='balanced_accuracy',return_train_score=
10     False)
11 grid.fit(X_training_data_PCA.to_numpy().astype('float')
12     ,np.transpose(y_training_data_PCA.to_numpy().
13     astype('int'))[0])
14 print("best mean cv score= ",grid.best_score_)
15 print("best parameters= ",grid.best_params_)
16 y_pred_knn=grid.best_estimator_.predict(X_test_pca)
17 print(confusion_matrix(y_test['HIV_active'],np.
18     transpose(np.matrix(y_pred_knn))))
19 X_test_pca=pd.DataFrame(data=pca.transform(X_test))
20 grid.score(X_test_pca,y_test['HIV_active'])
```

Listing 5: GridSearchCV-KNN-PCA [4]

```
1 best mean cv score= 0.7096760964856432
2 best parameters= {'metric': 'manhattan', 'n_neighbors': 1}
3 [[3757 80]
4  [ 74 70]]
5 0.7326307445052559
```

Listing 6: GridSearchCV-KNN-PCA results

You can see the results of the algorithm accuracy with different meters and hyperparameters in Figure 4.

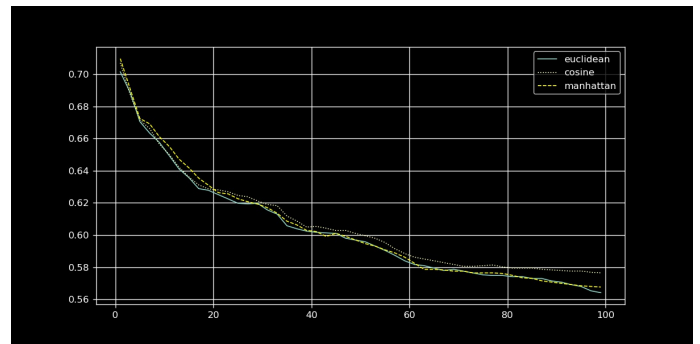


Figure 4: KNN accuracy

I implemented the KNN algorithm on the data without reducing the PCA dimensions, but it did not obtain the desired result, so I omitted the implementation of the algorithm on the data without PCA here. Also, in this project, I felt that PCA on some algorithms like KNN improves the result. While in the case of some algorithms, such as the decision tree algorithm, it can destroy the results. The reason for this is almost apparent given the nature of these two algorithms. Anyway, this is just a conjecture, and it does not mean that this is correct.

5.2 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. [9]

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. For example, I invite you to compare the results of Listings (7,8) and (9,10).

```
1 from sklearn.tree import DecisionTreeClassifier,
  export_graphviz
2 DecisionTree = DecisionTreeClassifier()
3 param_grid={'criterion':['gini','entropy'],'max_depth':
  :[18,20]}
4 grid=GridSearchCV(DecisionTree,param_grid=param_grid,
  cv=10,scoring='balanced_accuracy',
  return_train_score=False)
5 grid.fit(X_training_data_PCA.to_numpy().astype('float')
  ),np.transpose(y_training_data_PCA.to_numpy().
  astype('int'))[0])
6 print("best mean cv score= ",grid.best_score_)
7 print("best parameters= ",grid.best_params_)
8 y_pred_dt=grid.best_estimator_.predict(X_test_pca)
9 print(confusion_matrix(y_test['HIV_active'],np.
  transpose(np.matrix(y_pred_dt))))
10 print(grid.score(X_test_pca,y_test['HIV_active']))
```

Listing 7: GridSearchCV-Decision Tree-PCA [4]

```
1 best mean cv score= 0.6131812593469491
2 best parameters= {'criterion': 'entropy', 'max_depth': 20}
3 [[3725 112]
4  [ 100  44]]
5 0.6381830423073582
```

Listing 8: GridSearchCV-Decision Tree-PCA results

```
1 from sklearn.tree import DecisionTreeClassifier,
  export_graphviz
2 DecisionTree = DecisionTreeClassifier()
3 param_grid={'criterion':['gini','entropy'],'max_depth':
  :[19,18]}
4 grid=GridSearchCV(DecisionTree,param_grid=param_grid,
  cv=5,scoring='balanced_accuracy',
  return_train_score=False)
5 grid.fit(X_training_data.to_numpy().astype('float'),
  y_training_data['HIV_active'])
6 print("best mean cv score= ",grid.best_score_)
7 print("best parameters= ",grid.best_params_)
8 y_pred_dt2=grid.best_estimator_.predict(X_test)
9 print(confusion_matrix(y_test['HIV_active'],np.
  transpose(np.matrix(y_pred_dt2))))
10 print(grid.score(X_test,y_test['HIV_active']))
```

Listing 9: GridSearchCV-Decision Tree [4]

```
1 best mean cv score= 0.6578310426299441
2 best parameters= {'criterion': 'entropy', 'max_depth': 19}
3 [[3793  44]
4  [  92  52]]
5 0.6748219094778907
```

Listing 10: GridSearchCV-Decision Tree results

5.3 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, θ , where θ is sampled i.i.d.¹⁰ 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 11 and 12.

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.datasets import make_classification
3 RFC = RandomForestClassifier()
4 param_grid={'n_estimators':[200,300,1000],'criterion':
  :['gini','entropy'],'max_depth':[5,6,8,10,18,25],
  'max_features':['auto','sqrt','log2']}
5 grid=GridSearchCV(RFC,param_grid=param_grid,cv=5,
  scoring='balanced_accuracy',return_train_score=
  False)
6 grid.fit(X_training_data.to_numpy().astype('float'),
  y_training_data['HIV_active'])
7 print("best mean cv score= ",grid.best_score_)
8 print("best parameters= ",grid.best_params_)
9 y_pred_RFC=grid.best_estimator_.predict(X_test)
10 print(confusion_matrix(y_test['HIV_active'],np.
  transpose(np.matrix(y_pred_RFC))))
11 print(grid.score(X_test,y_test['HIV_active']))
```

Listing 11: GridSearchCV-Random Forest [4]

```
1 best mean cv score= 0.6348288515425576
2 best parameters= {'criterion': 'gini', 'max_depth': 25, '
  max_features': 'auto', 'n_estimators': 200}
3 [[3828  9]
4  [ 106  38]]
5 0.6307716532012857
```

Listing 12: GridSearchCV-Random Forest results

5.4 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. [10] In the following, I have implemented this beautiful algorithm on my data.

```
1 from xgboost import XGBClassifier
2 import warnings
3 warnings.simplefilter('ignore')
4 param_grid={'n_estimators':[700,2000],'max_depth':
  :[4,6],'eval_metric':['mlogloss']}
5 grid=GridSearchCV(XGBClassifier(),param_grid=
  param_grid,cv=10,scoring='balanced_accuracy',
  return_train_score=False)
6 grid.fit(X_training_data.to_numpy().astype('float'),(
  np.transpose((y_training_data['HIV_active']).
  to_numpy().astype('int'))))
7 print("best mean cv score= ",grid.best_score_)
8 print("best parameters= ",grid.best_params_)
9 y_pred_xgbc=grid.best_estimator_.predict(X_test)
10 print(confusion_matrix(y_test['HIV_active'],np.
  transpose(np.matrix(y_pred_xgbc))))
```

```
11 print(grid.score(X_test, y_test['HIV_active']))
```

Listing 13: GridSearchCV-XGBoost [4]

```
1 best mean cv score= 0.6785054406535646
2 best parameters= {'eval_metric': 'mlogloss', 'max_depth': 4,
3                   'n_estimators': 2000}
3 [[3817 20]
4  [ 90 54]]
5 0.684893797237425
```

Listing 14: GridSearchCV-XGBoost results

5.5 C-Support Vector Classification

In machine learning, support-vector machines (SVMs, also support-vector networks[11]) 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 param_grid={'kernel': 'poly', 'degree': [5, 6, 7, 8]}
3 grid=GridSearchCV(SVC(), param_grid=param_grid, cv=10,
4                   scoring='balanced_accuracy', return_train_score=
5                   False)
4 grid.fit(X_training_data.to_numpy().astype('float'), (
5   np.transpose((y_training_data['HIV_active']).
6   to_numpy().astype('int'))))
5 print("best mean cv score= ", grid.best_score_)
6 print("best parameters= ", grid.best_params_)
7 y_pred_svm=grid.best_estimator_.predict(X_test)
8 print(confusion_matrix(y_test['HIV_active'], np.
9   transpose(np.matrix(y_pred_svm))))
9 print(grid.score(X_test, y_test['HIV_active']))
```

Listing 15: GridSearchCV-SVC [4]

```
1 best mean cv score= 0.6878128245770883
2 best parameters= {'degree': 7, 'kernel': 'poly'}
3 [[3769 68]
4  [ 79 65]]
5 0.7168333550516897
```

Listing 16: GridSearchCV-SVC results

5.6 Logistic Regression

Logistic regression is the appropriate regression analysis to conduct when the dependent variable is dichotomous (binary). Like all regression analyses, the logistic regression is a predictive analysis. Logistic regression is used to describe data and to explain the relationship between one dependent binary variable and one or more nominal, ordinal, interval or ratio-level independent variables. [12]

In logistic regression we learn a family of functions h from \mathbb{R}^d to the interval $[0, 1]$. However, logistic regression is used for classification tasks: We can interpret $h(x)$ as the *probability* that the label of \mathbf{x} is 1. The hypothesis class associated with logistic regression is the composition of a sigmoid function $\phi_{\text{sig}} : \mathbb{R} \rightarrow [0, 1]$ over the class of linear functions L_d . In particular, the sigmoid function used in logistic regression is the *logistic function*, defined as

$$\phi_{\text{sig}}(z) = \frac{1}{1 + \exp(-z)}$$

The results of the Logistic Regression algorithm on my data were not very good. That's why I examined both the primary data and the data generated by the PCA algorithm with this algorithm. The result was better when I did not use PCA. But compared to other algorithms, it

was not very desirable. I would like to draw your attention to Listings 17, 18, 19 and 20.

```
1 from sklearn.linear_model import LogisticRegression
2 param_grid={'penalty': ['l1', 'l2', 'none', 'elasticnet'],
3             'solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', '
4             saga']}
3 grid=GridSearchCV(LogisticRegression(), param_grid=
4   param_grid, cv=5, scoring='balanced_accuracy',
5   return_train_score=False)
4 grid.fit(X_training_data_PCA.to_numpy().astype('float'
5   ), (np.transpose((y_training_data_PCA['HIV_active'
6   ]).to_numpy().astype('int'))))
5 print("best mean cv score= ", grid.best_score_)
6 print("best parameters= ", grid.best_params_)
7 y_pred_lr_p=grid.best_estimator_.predict(X_test_pca)
8 print(confusion_matrix(y_test['HIV_active'], np.
9   transpose(np.matrix(y_pred_lr_p))))
9 print(grid.score(X_test_pca, y_test['HIV_active']))
```

Listing 17: GridSearchCV-Logistic Regression-PCA [4]

```
1 best mean cv score= 0.5343843746116311
2 best parameters= {'penalty': 'none', 'solver': 'newton-cg'}
3 [[3835 2]
4  [ 130 14]]
5 0.5483504908348537
```

Listing 18: GridSearchCV-Logistic Regression-PCA results

```
1 from sklearn.linear_model import LogisticRegression
2 param_grid={'penalty': ['l1', 'l2', 'none', 'elasticnet'],
3             'solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', '
4             saga']}
3 grid=GridSearchCV(LogisticRegression(), param_grid=
4   param_grid, cv=5, scoring='balanced_accuracy',
5   return_train_score=False)
4 grid.fit(X_training_data.to_numpy().astype('float'), (
5   np.transpose((y_training_data['HIV_active']).
6   to_numpy().astype('int'))))
5 print("best mean cv score= ", grid.best_score_)
6 print("best parameters= ", grid.best_params_)
7 y_pred_lr=grid.best_estimator_.predict(X_test)
8 print(confusion_matrix(y_test['HIV_active'], np.
9   transpose(np.matrix(y_pred_lr))))
9 print(grid.score(X_test, y_test['HIV_active']))
```

Listing 19: GridSearchCV-Logistic Regression [4]

```
1 best mean cv score= 0.5862637041676992
2 best parameters= {'penalty': 'none', 'solver': 'lbfgs'}
3 [[3830 7]
4  [ 115 29]]
5 0.5997822734775432
```

Listing 20: GridSearchCV-Logistic Regression results

5.7 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 param_grid={'n_estimators': [500, 700, 1000, 8000], '
  algorithm': ['SAMME.R'], 'learning_rate': [1]}
3 grid=GridSearchCV(AdaBoostClassifier(), param_grid=
  param_grid, cv=5, scoring='balanced_accuracy',
  return_train_score=False)
4 grid.fit(X_training_data_PCA.to_numpy().astype('float'
  ), (np.transpose((y_training_data_PCA['HIV_active'
  ])).to_numpy().astype('int'))))
5 print("best mean cv score= ", grid.best_score_)
6 print("best parameters= ", grid.best_params_)
7 y_pred_adaboost_pca=grid.best_estimator_.predict(
  X_test_pca)
8 print(confusion_matrix(y_test['HIV_active'], np.
  transpose(np.matrix(y_pred_adaboost_pca))))
9 print(grid.score(X_test_pca, y_test['HIV_active']))
```

Listing 21: GridSearchCV-AdaBoost-PCA [4]

```
1 best mean cv score= 0.5711347496883823
2 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
  'n_estimators': 1000}
3 [[3811  26]
4  [ 123  21]]
5 0.5695286030753193
```

Listing 22: GridSearchCV-AdaBoost-PCA results

```
1 from sklearn.ensemble import AdaBoostClassifier
2 param_grid={'n_estimators': [500, 700, 1000, 8000], '
  algorithm': ['SAMME', 'SAMME.R'], 'learning_rate'
  : [1]}
3 grid=GridSearchCV(AdaBoostClassifier(), param_grid=
  param_grid, cv=5, scoring='balanced_accuracy',
  return_train_score=False)
4 grid.fit(X_training_data.to_numpy().astype('float'), (
  np.transpose((y_training_data['HIV_active'])).
  to_numpy().astype('int'))))
5 print("best mean cv score= ", grid.best_score_)
6 print("best parameters= ", grid.best_params_)
7 y_pred_adaboost=grid.best_estimator_.predict(X_test)
8 print(confusion_matrix(y_test['HIV_active'], np.
  transpose(np.matrix(y_pred_adaboost))))
9 print(grid.score(X_test, y_test['HIV_active']))
```

Listing 23: GridSearchCV-AdaBoost [4]

```
1 best mean cv score= 0.6508824113415643
2 best parameters= {'algorithm': 'SAMME.R', 'learning_rate': 1,
  'n_estimators': 8000}
3 [[3813  24]
4  [ 97  47]]
5 0.6600670011293546
```

Listing 24: GridSearchCV-AdaBoost results

5.8 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 overfitting¹¹. 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 param_grid={'n_estimators': [10, 100, 500, 5000], '
  max_samples': [1, 0.5], 'max_features': [1, 0.8, 0.6]}
```

```
3 grid=GridSearchCV(BaggingClassifier(), param_grid=
  param_grid, cv=5, scoring='balanced_accuracy',
  return_train_score=False)
4 grid.fit(X_training_data.to_numpy().astype('float'), (
  np.transpose((y_training_data['HIV_active'])).
  to_numpy().astype('int'))))
5 print("best mean cv score= ", grid.best_score_)
6 print("best parameters= ", grid.best_params_)
7 y_pred_bag=grid.best_estimator_.predict(X_test)
8 print(confusion_matrix(y_test['HIV_active'], np.
  transpose(np.matrix(y_pred_bag))))
9 print(grid.score(X_test, y_test['HIV_active']))
```

Listing 25: GridSearchCV-Bagging [4]

```
1 best mean cv score= 0.638850300575269
2 best parameters= {'max_features': 0.8, 'max_samples': 0.5, '
  n_estimators': 500}
3 [[3827  10]
4  [ 101  43]]
5 0.648002454174268
```

Listing 26: GridSearchCV-Bagging results

6 EVALUATION

In this section, I want to compare the accuracy of different algorithms 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 want to examine the final accuracy of the algorithm on the test set and cv score.

To select the final algorithm, it is obvious that we must select the best accuracy of the algorithm on the test set. Note Figure 5. As you can see, the KNN-PCA algorithm was the most accurate. The second place belongs to SVC.

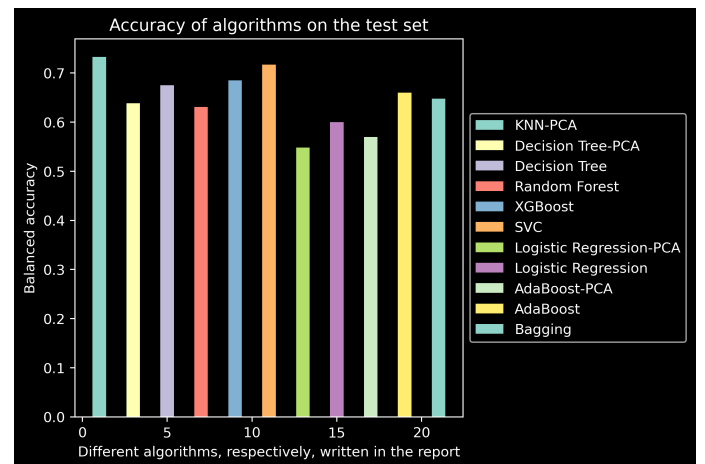


Figure 5: Accuracy of algorithms on the test set[15]

We have to keep in mind that the criteria for selecting hyperparameters were cv scores, so it would be better to show you a similar chart based on cv scores. Please pay attention to the Figure 6.

Looking at Figure 6, you may not notice the difference at first glance with Figure 5. For this reason, I have prepared Figure 7 for you, and all the bars in Figure 6 are displayed in a lighter gray color. All bars in Figure 5 are displayed in the initial format. You may notice a tiny

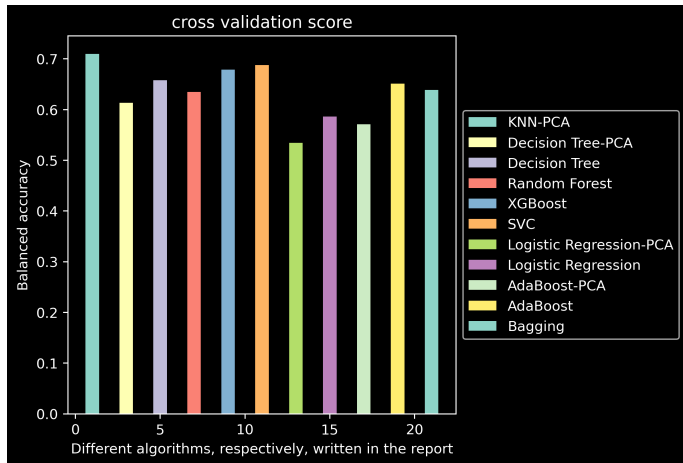


Figure 6: cross validation score[15]

amount of overfitting or even underfitting¹² in some algorithms. Of course, the amount is much smaller than what can be called overfitting or underfitting. The point that I find interesting or perhaps exciting is that the KNN algorithm with $n = 1$, which we have heard many times leads to overfitting, was surprisingly raised. Not only did overfitting not happen, but the cv score was even lower than the test score.

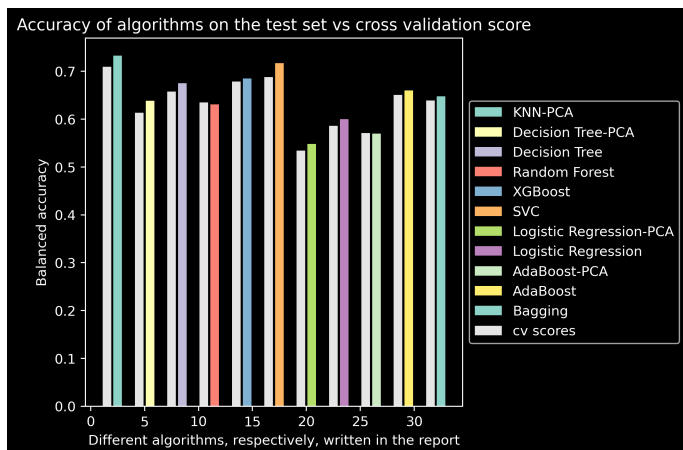


Figure 7: Accuracy of algorithms on the test set vs cross validation score[15]

A fundamental question between the two KNN-PCA and SVC algorithms, which were in the first and second place, respectively, which is really better? We can easily find the answer to this question by analyzing the confusion matrix.

$$\text{KNN-PCA : } \begin{bmatrix} \text{TN : 3757} & \text{FP : 80} \\ \text{FN : 74} & \text{TP : 70} \end{bmatrix} \quad \text{SVC : } \begin{bmatrix} 3769 & 68 \\ 79 & 65 \end{bmatrix}$$

The nature of the data is such that finding ones is more valuable than finding zeros. SVC performed better in finding zeros. In other words, it had less FP. Meanwhile, KNN-PCA performed about 7.5 percent better than SVC in finding ones. I think KNN-PCA is a better choice, given the data. However, I tried not to miss the chance of SVC as the second algorithm to be more accurate, but as the evidence shows, **KNN-PCA** is the winner.

Note: In the machine learning projects, principles oblige us to choose the best algorithm based on the validation set. This project has to select the best algorithm based on the Cv score. Here the best test set and Cv score are the same as the KNN-PCA algorithm. But basically, the best model should be based on the Cv score.

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

NOTES

¹This statistic is related to the year 2020.

²This statistic is related to the year 2020.

³This way, data labeled **zero** will not affect missing features labeled **one**.

⁴**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.[16]

⁵According to the group's agreement, the random number 1234 has been selected for `train_test_split`

⁶**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.[4]

⁷TP: True Positive, P: Positive, TN: True Negative, N: Negative

⁸FP: False Positive, FN: False Negative

⁹for example, when the training set is the entire Web and distances are based on links

¹⁰Independent and identically distributed

¹¹In 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.

¹²Underfitting is a scenario in data science where a data model is unable to capture the relationship between the input and output variables accurately, generating a high error rate on both the training set and unseen data.

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] unaids.org. Global hiv aids statistics, 2022.
- [3] hulab.rxnfinder.org. Convert smiles to image, 2019.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] Jr Mohammed J. Zaki, Wagner Meira. *Data Mining and Machine Learning-Fundamental Concepts and Algorithms*. Cambridge University Press, 2020.
- [8] Joydwp Mohajon. Confusion matrix for your multi-class machine learning model, 2020.
- [9] 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.
- [10] 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.
- [11] Corinna Cortes and Vladimir Naumovich Vapnik. Support-vector networks. *Machine Learning*, 20:273–297, 2004.
- [12] statisticssolutions.com. What is logistic regression?, 2022.
- [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] J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007.
- [16] 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.