# HEART DISEASE PREDICTION USING DATA MINING CLASSIFICATION ALGORITHMS

**S M RAHID HAQUE[1], MD. ATIK FOYSAL[1], ARUPKUMAR DAS[1], SHAHIDUL ISLAM LEON[1], DR. MD. ABDULLAH - AL - JUBAIR[2].**

[1]STUDENT OF COMPUTER SCIENCE DEPARTMENT,
[2]ASSISTANT PROFESSOR OF COMPUTER SCIENCE DEPARTMENT,
AMERICAN INTERNATIONAL UNIVERSITY-BANGLADESH, DHAKA.

**Abstract—** a range of conditions which affect heart is called heart diseases or "cardiovascular diseases". This disease can bring out heart attack, chest pain, stroke etc. By reviewing some research paper related to heart disease prediction it was identified that most of the paper using singular algorithm to predict the disease using machine learning algorithm. Some of them indicates that they can't use optimization techniques to improve their model performance. For these results, they have faced some problem to predict heart disease in an efficient manner by using their proposed system. To overcome these problems and for getting more accurate results in this medical study is very crucial that's why four different classification algorithms were implemented to predict heart disease and find out the effectiveness of these algorithms. In this study principal component analysis (pca) dimensionality reduction technique was applied which helped to get better results with the aim of better accuracy by using these algorithms since medical diagnosis is sensitive. For this approach data was collected from uci repository which was found in kaggle and it is named as "heart disease uci". It was observed that without principal component analysis (pca) logistic regression performed best to predict heart disease with principal component analysis (pca) k-nearest neighbors (knn) achieved greater accuracy compared to other classification algorithms. By applying pca it was identified that accuracy for other algorithms like decision tree and naive bayes also increased compared to originally.

**Index terms**— cardiovascular diseases, machine learning, pca, heart disease prediction, classification algorithm

## I. INTRODUCTION

Proper health is the key to a country's well-being. A range of conditions that affect the heart are called heart diseases. It is also called "cardiovascular diseases". In this disease, blood vessels can be prevented from moving and it can bring about heart attack, chest pain, stroke etc. Some other conditions also considered as heart diseases such as the heart state which attack the heart's muscle valves or pattern.

The concept of 'data mining' describes the process of extracting patterns from vast quantities of data. Data mining includes variety of algorithms like classification for supervised learning and clustering for unsupervised learning.

### A. Problem statement

In previous clinical trials, machine learning approaches were used to forecast heart disease. However, these researches concentrated on the singular results of specific machine learning techniques but not using any optimization of techniques for improvement [1] [2]. In addition, these researchers tried to use hybrid optimization methods for an optimized classification of machine learning. For heart disease prediction using data mining techniques some of the researches don't provide valid accuracy results [3] as they haven't provided any

---

This paragraph of the first footnote will contain the date on which you submitted your paper for review. It will also contain support information, including sponsor and financial support acknowledgment. For example, "This work was supported in part by the U.S. Department of Commerce under Grant BS123456."

The next few paragraphs should contain the authors' current affiliations, including current address and e-mail. For example, F. A. Author is with the National Institute of Standards and Technology, Boulder, CO 80305 USA (e-mail: author@ boulder.nist.gov).

S. B. Author, Jr., was with Rice University, Houston, TX 77005 USA. He is now with the Department of Physics, Colorado State University, Fort Collins, CO 80523 USA (e-mail: author@lamar.colostate.edu).

T. C. Author is with the Electrical Engineering Department, University of Colorado, Boulder, CO 80309 USA, on leave from the National Research Institute for Metals, Tsukuba, Japan (e-mail: author@nrim.go.jp).

optimization techniques in improving the classification model and haven't checked the performance of the model in the obtained result. So, there is a need for study to find out effectiveness and how classification algorithms work in predicting heart disease.

B. **Objectives**

The overarching goal of this study is to see whether machine learning classification algorithms will accurately predict heart disease. Based on the reviews of previous research the objectives of this research are to predict heart diseases using data mining classification techniques such as naive bias, k-nearest neighbor, logistic regression decision tree. Then online open-source platform is used which is https://jupyter.org/ where it can develop the proposed machine learning algorithms. Finally, it must figure out how to improve the model, compare the accuracy and performance of different classification models.

## II. BACKGROUND STUDY

According to world health organization, heart disease has been one of the top killers, with millions of individuals dying annually. This is a type of disease that attacks both the heart and vascular system of the body and causes a range of heart related infections like coronary heart disease, hypertension, cerebrovascular disease etc. [4]. According to earlier findings, the chances of developing coronary heart disease was 34.9 percentage for males and 24.2 percentage for females at the age of 70 and 42.4 at age 40 respectively [5]. Heart disease shows early signs of symptoms like breathlessness, orthopnea, fatigue, tiredness, ankle inflammation etc. And diagnosis includes chest radiograph, electrocardiogram (ecg), echocardiography etc. Which is really pricey [6].

In this modern age there is data in nearly all regards of our life like it can be in hospitals, business stores, social media where large volumes of data are collected like medical records, sales records etc. And data mining can help turn those data into a knowledge [7]. Data mining is seen as an adaptation in information technology from early database systems in 1970 where data are tabulated in relational database analysis systems [8]. Then in 1980 it was researched and progressed to advanced database storage techniques where there was an abundance of storage systems and equipment to store large amounts of data and finally in 1990 web-based databases like xml was introduced where data was plenty and information gathering using data mining now far easier. Data mining classification algorithms can be used to predict the labeled class of any data set and to solve real-world problems including heart disease [7] [8].

## III. METHODOLOGY

In the study, heart diseases prediction framework is introduced and this approach aims to analyze performance of different classification algorithms and find out effectiveness of the classification algorithms.

Datasets are encrypted and loaded within this proposed system. After that the data is divided into two sections, with one half of the sample being used for training and another half being used for evaluating our model. Then different classification techniques are used such as decision tree, k-nearest neighbor, logistic regression and naive bayes analyze performance by checking accuracy, precision, recall, f1 score, confusion matrix diagram of the algorithms to find the most effective classification algorithm. At last, pca feature selection technique is applied to improve performance of the algorithm.
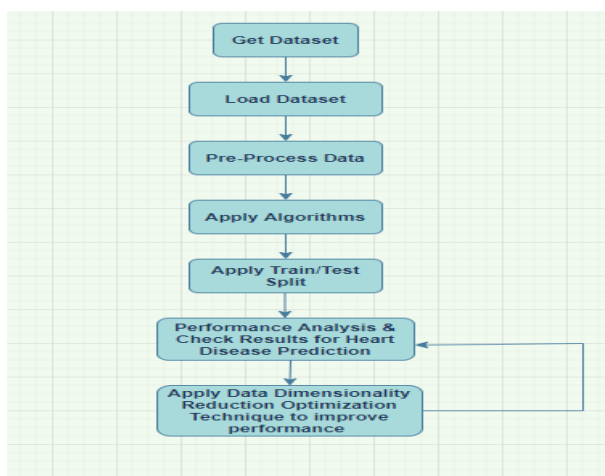
fig. 1. Proposed framework of heart disease prediction

A. Dataset collection

For this research, collection of data is the first phase of the proposed technique to predict the heart disease. In this process, the data will be collected from uci repository which was found in kaggle and it is named as "heart disease uci" [9]. This public data is widely used by ml researchers around the world for heart disease research in data mining as stated in the website. The "target" column holds if a patient has heart disease or not and it denoted by 1 or 0 which will be used for prediction.

This data has 13 attributes (age, gender, cp, trest bps, cholesterol, fbs, rest ecg, thalach, exang, old peak, slope, ca, thal) 1 class label attributes (diagnosis of heart disease) and 303 instances of data originally.

B. Loading and pre-processing data

After collecting the data second phase is about loading the data by reading the .csv file. For reading the .csv file the command used for it was the pd.read and pd was imported as import pandas as pd. Then the data was ultimately loaded inside the notebook for pre-processing.

After loading the data next step is about pre-processing it. The author in [11] said that in actuality datasets are extremely vulnerable because of noise, absent and conflicting data due to natural error or large datasets. Data pre-processing is the way to solve these issues. In this phase there are many preprocessing approaches that can used in the dataset. It can be data cleaning, data reduction, data integration etc. In visualization we have to check if the target data is balanced or not. In this research we are tracking the dataset to detect any null values. In order to do that isna() function will be used to detect null and in addition the sum of null values with .sum() will be calculated for each columns. Moreover, we will also apply the feature scaling method. This method is mainly used to normalize the heart disease dataset. In order to that library must transferred from sklearn which is standardscaler and this function will be loaded inside a variable so that the object can be created to apply this method.

C. Performing train-test split on algorithms

After pre-processing train-test split will be applied in the dataset. As stated by [10], it is necessary to train the data based on training data and make a prediction on the test data. So, the test data is basically used to assess the performance of the model.

In this research, dataset is divided into 70 percentage for training and remaining 30 percentage for testing purpose. Therefore, among the 303 instances roughly 91 samples will be used for testing purposes and other samples for training. Random state is also specified which is 40 for data splitting so that model perform with repeated data every time we run. To accomplish all of this, we must first install the scikit-learn library and then access the train test split() function from the python library.

In order to do perform the train-test split next step is to find the suitable algorithms. We have to apply train-test split on wide variety of supervised data mining classification algorithms. Classification algorithms that will be used are decision tree, knn, logistic regression and naive bayes. For all the algorithms again python machine learning library will be used and transfer the suitable algorithms.

### D. Checking performance and results of algorithms

In the second last phase, after applying algorithms next step will be is to check the performance of each of them. In this phase, accuracy score must be determined first. The number of correctly estimated instances in comparison to the overall number of instances will be determined by the accuracy ranking. Then confusion matrix will be applied to check the performance of each algorithm in more detail. Precision, recall, f1 score, and accuracy will be created by the confusion matrix, and the results will be mentioned in the results section of the confusion matrix article.

### E. Applying pca dimensional reduction technique

In the final phase, pca optimization techniques will be used and compare the result again with applying it. Pca is used to increase the algorithm's efficiency by making it work quicker by dramatically reducing training time, and it does so by reducing measurements [12]. It also helps reducing overfitting so algorithms like decision tree with low number of features in dataset can work properly with improved accuracy and performance.

## IV. IMPLEMENTATION AND EVALUATION

### A. Algorithms

#### i. Decision tree

In order to implement the decision tree certain steps are followed in our study. At first the decision tree is loaded using scikit-learn libraries and other required libraries. Then the feature and target variables are separated and train-test split is also applied. In order to build decision tree performance is improved by setting parameters with pre-pruning by selecting the max depth and attribute selection measurement like entropy [13]. Then the performance is evaluated using confusion matrix.

#### ii. Logistic regression

The primary goal of this technique is to estimate the probabilities of events. When the address variable is categorical and only has a few distinct variations, this method is used. In order to apply in this following research at first the logistic regression libraries are from scikit-learn from python. Then the data is loaded and split using train-test method. Then finally confusion matrix will be generated to determine the performance of the model.

#### Iii. K-nearest neighbor

In order to implement knn in our study python scikit-learn technique was used. At first the data was loaded then train/test spilt was applied. Knn was imported from the library and fitted into the data. One of the crucial steps in applying knn is to determine the value of 'k' [14]. In order to find the value of 'k' for loop was applied to test the accuracy in range of neighbors from 1 to 100 in our research. K was determined by plotting the graph with the respect to value of k and testing accuracy. The value of 'k' was chosen where accuracy was highest with small value of k to minimize overfitting.

#### Iv. Naïve bayes

In this study dataset composed of binary labels which is "1" or "0" so it is best to use gaussian naïve bayes technique to apply in the dataset to get the heart disease prediction. It also only measures the mean and standard deviation from the training dataset and it is easier to apply. So, in order to do that we used sklearn techniques to import gaussian naïve bayes from python library. Train/test method was applied and algorithm is fit inside the dataset. Finally, the prediction was made and the accuracy value was assessed.

### B. Performance measurement using confusion matrix

In order to get more information about the performance of the algorithms on accuracy is confusion matrix and it summarizes overall performance of the classification algorithms that is used in this research.

1) true positive (tp): true positive occurs when actual value was positive and predicted label value was also positive. In this case it means a patient has heart disease and algorithm also predicted that he/she actually has it.

2) true negative (tn): true negative represents when actual value was negative and predicted label was also negative. For instance, in this case, it means patient does not have any heart disease and also predicted to not have any.

3) false positive (fp): false positive happens when actual value was negative but the algorithm prediction came positive. For example, in this research it can be said that patient do not have any heart disease but predicted to have the disease

4) false negative (fn): false negative happens when actual value was positive but the algorithm prediction came negative. In this research, it can be said that patient have heart disease in their body but predicted to not have any.

5) accuracy: one of the common performance detection techniques used by confusion matrix is the accuracy measurement. It determines how frequently the algorithm can yield actual output. It can be measured as the ratio of number of correction prediction to the total number of predictions that is obtained by the classifier.

Accuracy: (t p + t n)/(t p + t n + f p + f n)  (1)

6) precision: precision measures how much percentage of patients that actually predicted to have heart disease among the patients who are classified to have heart disease. It can be calculated as:

precision: t p/(t p + f p)  (2)

7) recall: recall calculates how much patients have predicted to get heart disease among the patients that actually have this disease. It can be calculated as:

Recall: t p/(t p + f n)  (3)

8) f1-score: the f1-score is taken as the harmonical mean of accuracy as well as recall. The maximum value is given when the same accuracy and reminder value are the same.

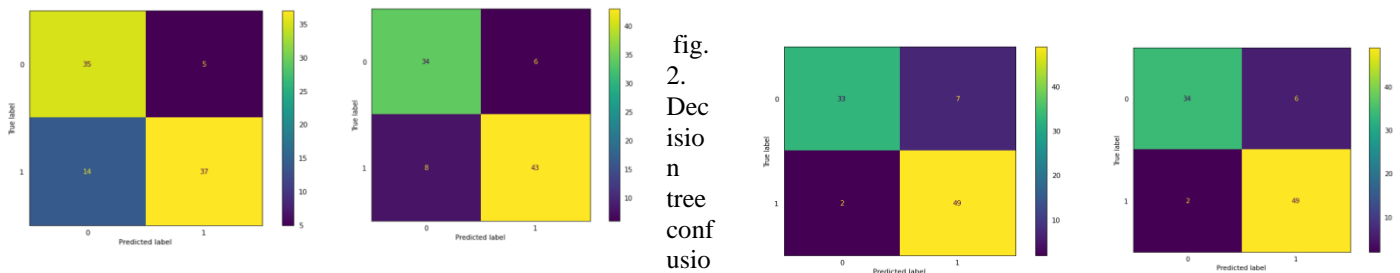F1−score: (precision ∗recall)/(precision+recall) ∗ 2  (4)

## C. Train-test split

In order to implement train/test split in the research we used scikit-learn technique from python library by passing the train test split() function and applied it to solve our classification problems. In this study, we also used 70 percentage for training the data and other 30 percentage for testing purpose and also 42 selected as random state so that same data used every time instead of using different.

## D. Result evaluation

1) performance of decision tree before and after applying pca:

The figure below demonstrates the original results of decision tree(left) and result after applying decision tree with pca (right). It can be interpreted that 35 patients have heart disease positive which is true positive with 14 patients false negative originally without pca whereas 34 predicted to have heart disease correctly and only 8 patients with false negative after pca which is an improvement. Therefore, the achieved recall value is increased which is 0.714 to 0.809. Decision tree also performed relatively well compared to original as the accuracy is increased from 0.791 to 0.846 which is quite significant. The total result of recall and f1-score was also shown to have a beneficial influence with pca.



fig. 2. Decision tree confusion matrix before (left) & after applying pca (right)

2) performance of logistic regression before and after pca:

The figure below illustrates the initial result of logistic regression (left) and result of this algorithm with pca (right). It can be seen that 35 patients predicted to have heart disease but only 2 patients are incorrectly predicted to not have the disease without pca but 36 patients predicted whereas again 5 patients wrongly predicted in terms

of true positive and false positive with pca. It can be seen that there is a change in accuracy as it is decreased from 0.923 to 0.901. So, overall, there is a decrease in progress after applying pca therefore it is not recommended to use pca with logistic regression as it performs relatively worse.
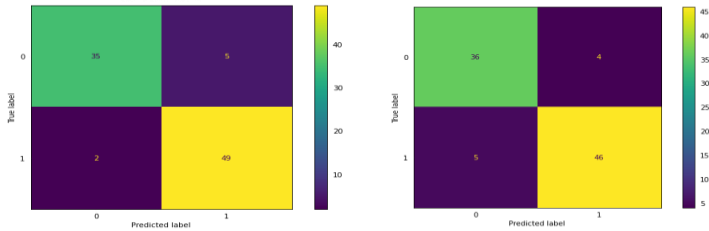


Fig. 3. Logistic regression confusion matrix before (left) & after applying pca (right)

3) performance of knn before and after pca:

In the results below, it depicts the results of confusion matrix of knn algorithm before(left) and applying (right) pca technique. Initially, it can be seen that without pca 34 patients have predicted to have heart disease but only 2 incorrectly identified to don't have heart disease and 7 patients incorrectly predicted as heart disease patient. But in case of pca, 6 patients wrongly predicted to have this disease, so pca performed slightly better compared to initially. The initial result also shows that knn performed with the accuracy of 0.901 but with applying pca it is increased to 0.912. Recall, precision and f1-score also increased slightly with pca so it is better to use pca with knn algorithm as it performs relatively better.

Fig. 4. Knn confusion matrix before (left) & after applying pca (right)

4) performance of naïve bayes before and after pca:

The figure below represents the results of confusion matrix before and after applying pca. It can be seen that in both cases 5 patients wrongly predicted to have heart disease. However, without pca, six patients were incorrectly predicted to be safe of heart disease, while only five patients were incorrectly predicted to be free of heart disease with pca. Accuracy of the result in predicting heart disease also slightly increased after applying pca from 0.879 to 0.890. Recall, precision, f1 score all increased drastically. So, overall, it can be seen that there is an improvement in the result after applying pca and no result seen any decrease in value after applying.
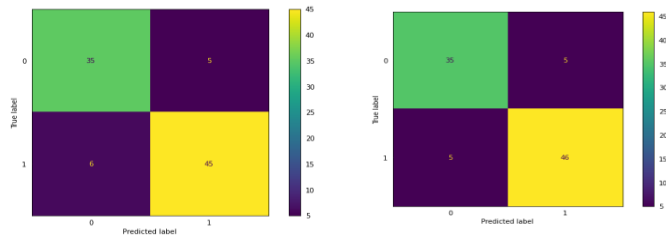


Fig. 5. Naïve bayes confusion matrix before (left) & after applying pca (right)

E. Findings

After applying dimensionality reduction technique with pca there are 9 attributes remaining in the dataset based on the explained variance. In the dataset those 9 attributes are trest bps, cholesterol, fbs, rest ecg, thalach, exang, old peak, ca, thal with target attributes. From the above result evaluation proof is provided that logistic regression outperformed other models in heart disease diagnosis, with an accuracy of 0.923. Other algorithms like knn and naïve bayes also performed well and results are close to logistic regression with accuracy 0.901 and 0.879 respectively. But decision tree performed worse out of all the four algorithms with only 0.791 accuracy which is very low compared to others. After the development of pca with these algorithms it can be seen that knn

performed well as accuracy increased to 0.912 compared to other algorithms. One of the highlights of using pca is that decision tree seen huge improvement in terms of accuracy as it increased from 0.791 to 0.846 so it can be concluded that reducing dimension of the dataset with pca can increase performance of decision tree significantly. Naïve bayes also seen slight improvement while using pca with accuracy of 0.890. In introducing pca performance metrics for the logistic regression decreased all over the board. But there is an improvement for all the other three algorithms which is decision tree, knn and naïve bayes by applying pca. As dimensionality of the dataset reduced with pca so the accuracy of the decision tree increased greatly since number of leaf reduced. But with introduction of pca logistic regression performance reduced, since minimizing the number of features in the dataset it resulted in reducing our accuracy slightly. By applying pca for knn and naïve bayes we can see that there are slight changes in accuracy. The table below shows the results of the accuracy before and after applying pca:

| Algorithms | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|
| Decision tree | 0.846 | 0.850 | 0.809 | 0.828 |
| Logistic regression | 0.901 | 0.900 | 0.878 | 0.889 |
| K-nearest neighbor | 0.912 | 0.850 | 0.944 | 0.895 |
| Naïve bayes | 0.890 | 0.875 | 0.875 | 0.875 |

TABLE I
RESULT OF ACCURACY, PRECISION, RECALL, F1-SCORE WITHOUT PCA

| Algorithms | Accuracy | Precision | Recall | f1-score |
|---|---|---|---|---|
| Decision tree | 0.791 | 0.875 | 0.714 | 0.786 |
| Logistic regression | 0.923 | 0.875 | 0.945 | 0.908 |
| K-nearest neighbor | 0.901 | 0.825 | 0.943 | 0.880 |
| Naïve bayes | 0.879 | 0.875 | 0.853 | 0.864 |

TABLE II
RESULT OF ACCURACY, PRECISION, RECALL, F1-SCORE WITH PCA

## V. LIMITATION AND FUTURE DIRECTIONS

In this research we are limited to use aged dataset from kaggle and the dataset size was small too. It contains only 303 features of data. We are separated from each other due to pandemic so it is harder to collect any recent data together for us to apply. It is better to use large data large data for prediction as it is difficult to know how quick the algorithm runs with small data and compare with pca and each algorithm.

In future we want to collect data from cardiac hospitals from bangladesh when the pandemic ends and apply the algorithms to see the results for prediction. It will help to analyze the effectiveness of the results in our country's perspective. We also want to demonstrate the algorithms' success in a roc curve mode. It will help to analyze the results in graphical forms which will help to easily identify results from the curves. We are also keen on developing a hybrid model by using the mixture of some of the algorithms that we used in this research to get better and accurate heart disease prediction. There are also more feature selection optimization techniques that we can apply despite only using pca in future. We proposed to apply mutual information gain feature selection, chi-square method, correlation heat map etc. In future to find more precise results of the algorithms that we applied. Also, we can apply more machine learning algorithms like svm, random forest and gradient boosting algorithms for more results for comparison in future works.

## VI. CONCLUSION

In this study there were different classification algorithms which were used to predict heart disease and analyzed performance of these algorithms. Pca was also optimized as an optimization technique to get better result for this strategy. By analyzing all the classification algorithms, it can be interpreted that logistic regression performed very well compared to other classification techniques. But with pca all results of accuracy like decision tree, knn and naïve bayes seen better results compared to originally.

[1] REFERENCES

[2] Khourdifi, y., bahaj, m. (2019). Heart disease prediction and classification using machine learning algorithms optimized by particle swarm optimization and ant colony optimization. International journal of intelligent engineering and systems, 12(1), 242–252. Https://doi.org/10.22266/ijies2019.0228.24 ding, w. And marchionini, g. 1997 a study on video browsing strategies. Technical report. University of maryland at college park.

[3] Patel, j., tejalupadhyay, s., patel, s. B. (2016). Heart disease prediction using machine learning and data mining technique. Journal - ijcscvolume: 7.

[4] Rajesh, n., t, m., hafeez, s., krishna, h. (2018). Prediction of heart disease using machine learning algorithms. International journal of engineering technology, 7(2.32), 363.

[5] South eastern health and social care trust. (2013). Cardiovascular disease. Communications department

[6] Lloyd-jones, d. M., larson, m. G., beiser, a., levy, d. (1999). Lifetime risk of developing coronary heart disease. The lancet, 353(9147), 89–92.

[7] Ministry of health kenya. (2015). Kenya national guideline for cardiovascular diseases management. Division of noncommunicable diseases.

[8] Bramer, m. (2020). Principles of data mining (undergraduate topics in computer science) (4th ed. 2020 ed.). Springer.

[9] Han, j., kamber, m., pei, j. (2011). Data mining: concepts and techniques (the morgan kaufmann series in data management systems) (3rd ed.). Morgan kaufmann.

[10] Heart disease uci. (2018, june 25). Kaggle. Https://www.kaggle.com/ronitf/heart-disease-uci

[11] Bronshtein, a. (2020, march 24). Train/test split and cross validation in python - towards data science. Medium. Https://towardsdatascience.com/train-test-split-and-cross-validationin-python-80b61beca4b6

[12] Acuna, e. (2011). Preprocessing in data mining. International encyclo- ˜ pedia of statistical sci-ence, 1083–1085.

[13] Kumar, n., kumar, n., profile, v. M. C. (2019). Advantages and disadvantages of princi-pal component analysis in machine learning. The professionals point.

[14] Navlani, a. (2018, december 28). Decision tree classification in python. Datacamp

[15] Sanjay.m (2018, november 2). Machine learning — knn using scikit-learn - towards data science. Medium. Https://towardsdatascience.com/knn-using-scikit-learn-c6bed765be75