

PREDICTING THE ONSET OF DIABETES BASED ON DIAGNOSTIC MEASURES

FYS-STK3155: PROJECT 3

Kim-Roger Grønmo
<http://github.com/kimgronmo/>

December 16, 2020

Abstract

In this paper we use data from the Pima Indians Diabetes Database to predict the onset of diabetes based on diagnostic measures. This is a binary classification problem which is solved by using different methods of classification as provided by the standard libraries in Scikit-learn. We find that the Gradient Boosting Classifier is best suited with an accuracy score of 0.90, while XGBoost and the soft Voting Classifier follows close behind with an accuracy score of 0.89. It is worth noting that the dataset had some issues with missing values. How this is resolved by either replacing or removing samples can impact the performance of the evaluated models and is worthy of further investigations.

1 Introduction

Our modern lifestyle comes with certain costs. A more sedentary lifestyle with less exercise and easier access to high calorie content foods, causes increases in lifestyle diseases. A common one that is increasing in recent years is diabetes. This disease has complications which include health problems such as kidney disease, nerve damage, retinal disease, heart disease and stroke.

This provides a strong motivation for quickly identifying and preventing this disease. Although doctors have methods available for identifying diabetes in a patient, a big part of the problem is for the patients themselves

to be aware of risk factors and what they themselves can do to prevent the onset of the disease. Machine learning algorithms might in the future enable doctors to clearly show patient how changes in a certain predictor variable, such as weight loss or lower blood sugar levels, can change the probability of either getting the disease or recovering from it.

Although the use of machine learning algorithms in practical treatment and disease preventions has its share of controversies, it is an important topic for future research.

Our goal in this project is to build and evaluate machine learning models in order to accurately predict whether or not the patients in the dataset have diabetes or not. As there are a lot of models to choose from we will not go into great depth of how each of them functions, but instead give a quick overview of each of them and the methods used to evaluate them. For further information we refer the reader to the standard Scikit-learn documentation and the references provided in our reference section.

We begin this report by looking at some theory behind the models and methods we will use, including a brief overview of a decision tree. We then investigate our data set and do some data cleaning and feature engineering. We present some of our findings and discuss the results. The report is then concluded with a summary of our findings.

2 Theory and Methods

The general references for this section is reference [1], [2] and [3]. For discussions about how each of these methods are implemented in Scikit-learn we refer the reader to the tool-kits homepage at [4]. We give a brief overview of each of the methods.

2.1 Decision Trees

2.2 Evalution Methods for Classification problems

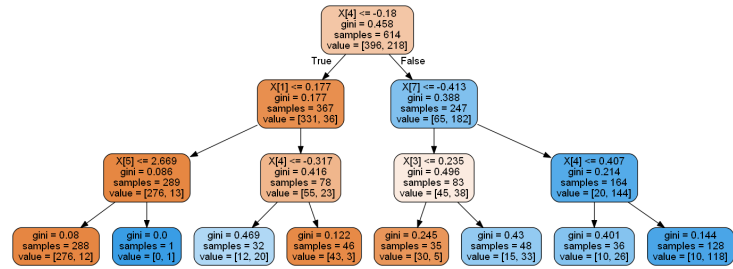
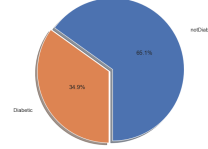
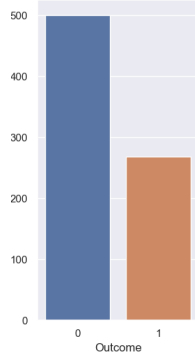


Figure 1: Example tree from the diabetes database with depth 3

2.3

2.4



(b) Distribution of patients

(a) Number of patients with diabetes (1) and non diabetic (0)

Figure 2: Overview of the patients in the data set

3 The Pima Indians Diabetes Dataset

This dataset is from the National Institute of Diabetes and Digestive and Kidney Diseases. All patients in this set are females at least 21 years old of Pima Indian heritage. The objective of the set is to diagnostically predict whether or not a patient has diabetes.

The dataset is available at [5]. A copy of it is also in the DataFiles directory in the github project folder.

3.1 Data Exploration and feature engineering

The datasets consist of several medical predictor variables and one target variable, Outcome. Outcome 1 represents the patient having diabetes, while 0 is non diabetic. The feature variables include the number of pregnancies the patient has had, their BMI, insulin level, age, and so on. Note the diabetes pedigree function is a function which scores likelihood of diabetes based on family history.

Figure 3 shows the distribution of patients among the two outcomes. We can see that the data set is not very imbalanced. No single outcome greatly outnumbers the other. Figure 4 shows the correlation between different features. We can see that Glucose levels have the biggest influence on the outcome, followed by insulin levels and BMI. From further exploration of the data set some obvious problems became apparent.

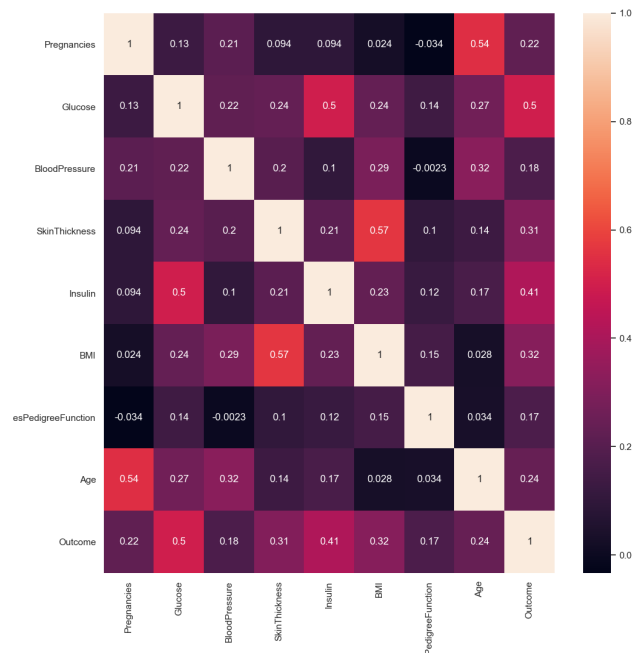


Figure 3: Correlation of different features

Although there were no missing values, certain factors were set to 0 indicating missing values. Features such as skin thickness, bmi etc are not possible to be at 0 for a living person. This was dealt with by splitting the data set in two (diabetic vs non diabetic), and replacing the 0 values with the average of the respective column for the given group. The two parts were then merged back into a complete set. Since the data set was rather small this seemed a better compromise than dropping the samples with the missing values.

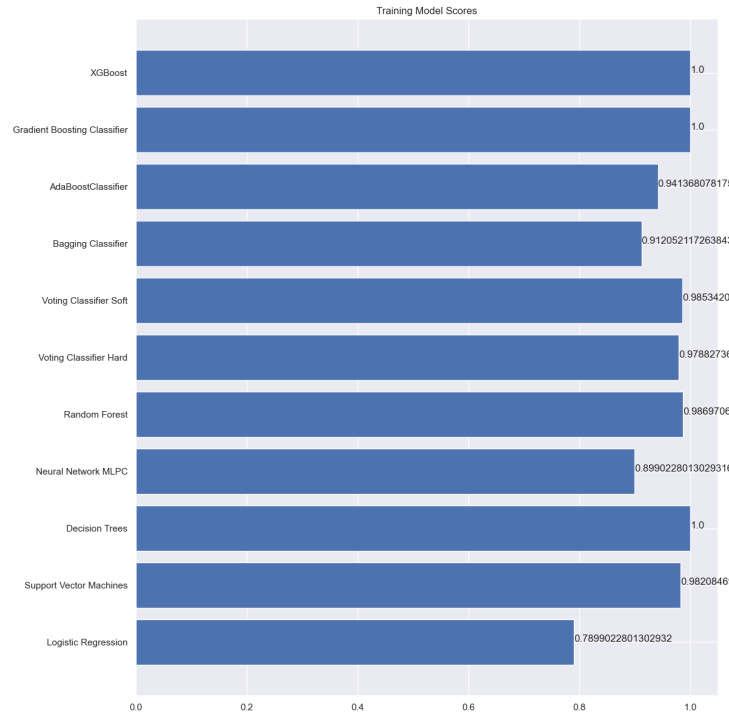


Figure 4: Training data model accuracy scores

4 Results and Discussion

Further results for running the code are included in the github folder. Figure 4 and 5 give the accuracy scores for the training data and test data. We can see that several models reach 100 percent accuracy with the training data. This accuracy for the test data is not as good so the models are overfitting on the training data. This is apparent for Gradient Boosting, XGBoost and Decision Trees.

We can see that the best accuracy score for the test set is given by Gradient Boosting at 0.902, while Decision Trees, XGBoost and Voting Classifier Soft are right behind with 0.89. In order to differentiate these models we will look closer at the confusion matrices to see how well they can predict the actual disease and not just the cases where there are no disease outcome.

From Figure 6 and 7 we can see there are clear differences in how well the models perform in identifying critical outcome (outcome 1: diabetes).

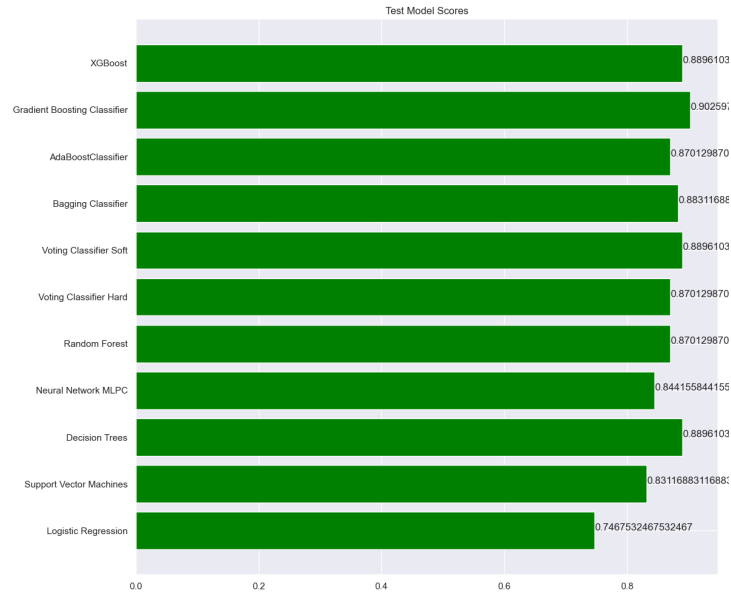
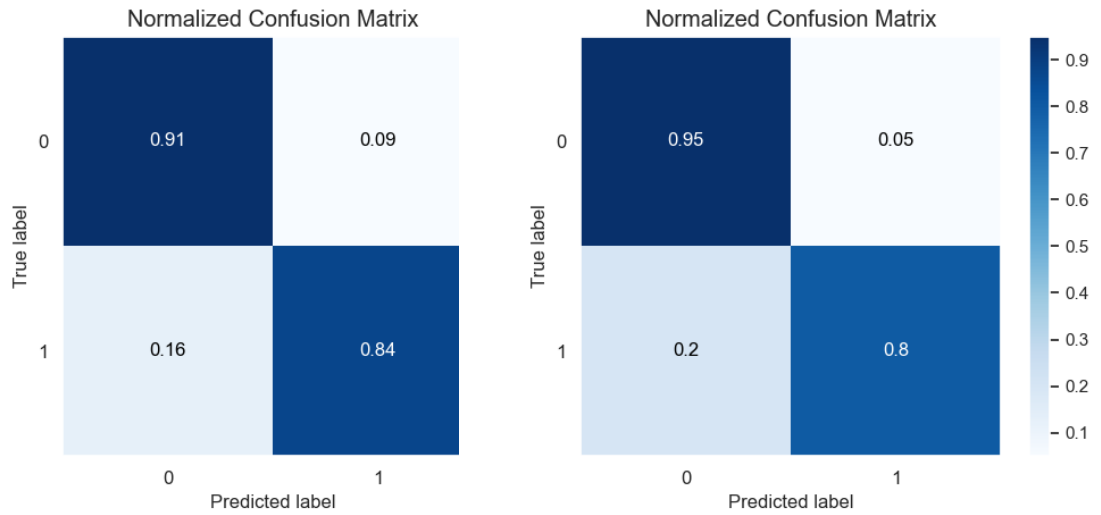


Figure 5: Test data model accuracy scores



(a) Voting Classifier Soft

(b) Gradient Boosting Classifier

Figure 6: Confusion Matrices

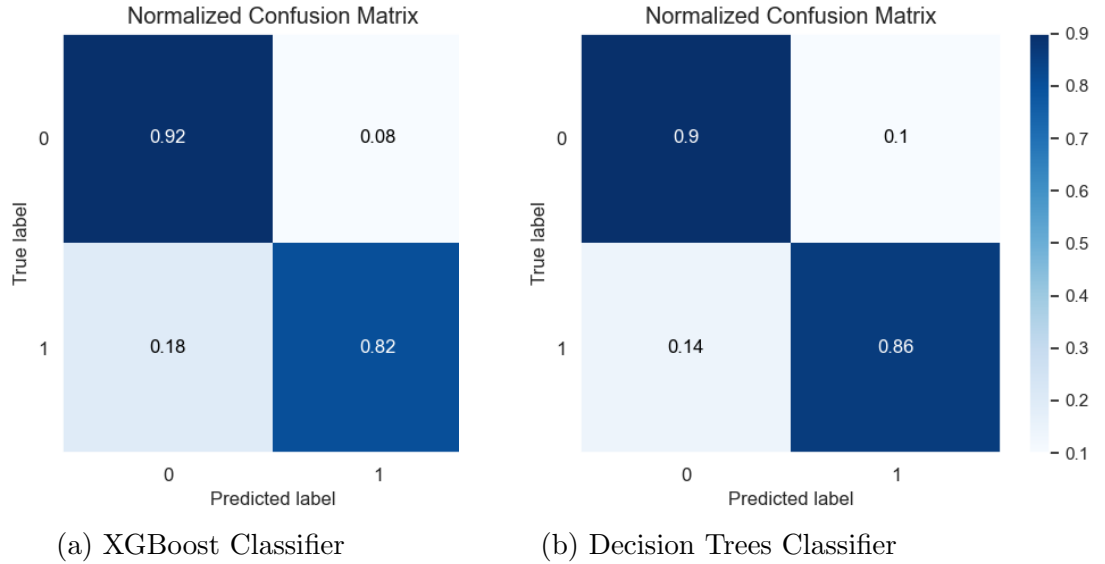


Figure 7: Confusion Matrices

5 Conclusion

We have evaluated different models for a binary classification problem. We found that although several of these models had almost the same accuracy score, the DecisionTree Classifier was slightly better at correctly identifying diabetes outcome. In cases of illness it is often important to quickly administer treatments. It is better to be misclassified as a diabetic case and not have the disease, than it is to be classified as non diabetic and have the disease.

The dataset available had some issues with missing data. Combined with a somewhat small sample size this can effect the accuracy of our models. In order to improve our models we might investigate further how to split the data differently and also how to replace the missing values by for instance using an average of their nearest neighbours, rather than sample averages. Further investigation into this matter is necessary to draw a proper conclusion.

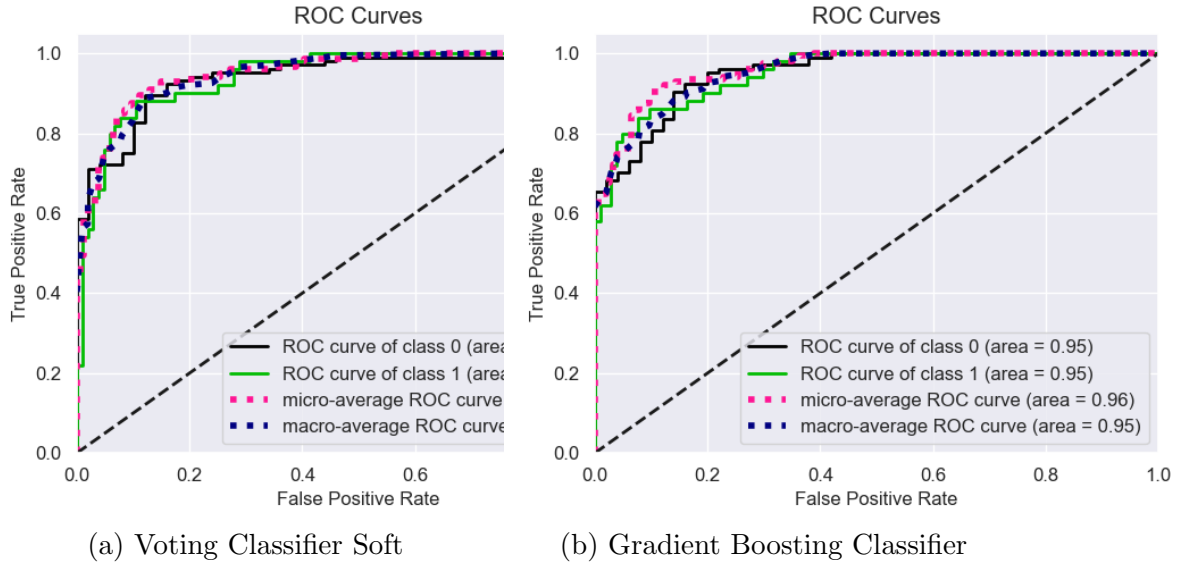


Figure 8: Receiver Operating Characteristic Curve

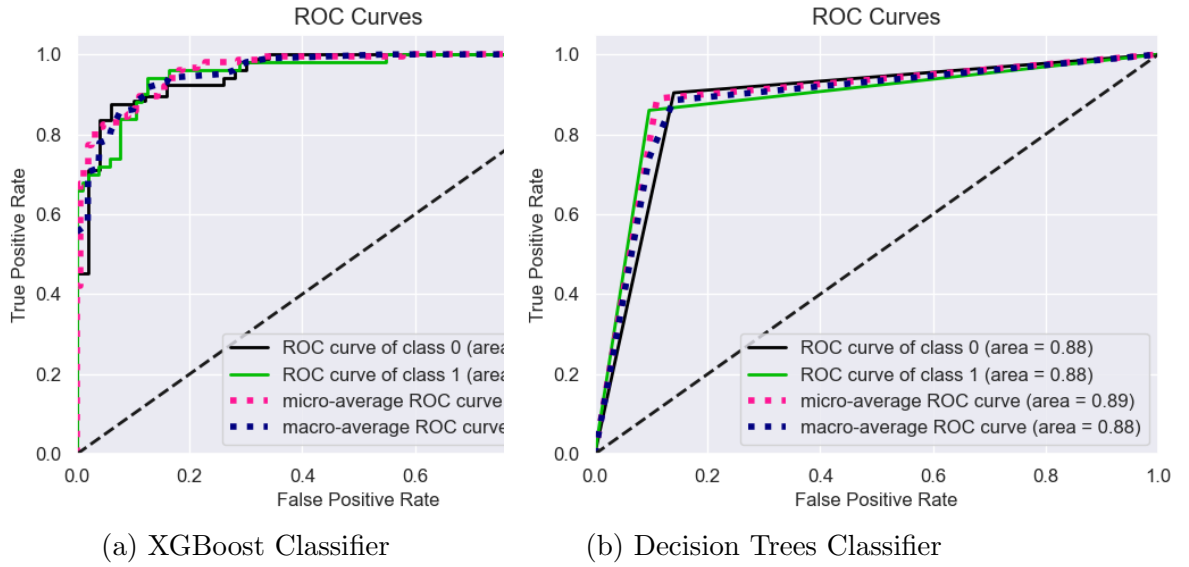


Figure 9: Receiver Operating Characteristic Curve

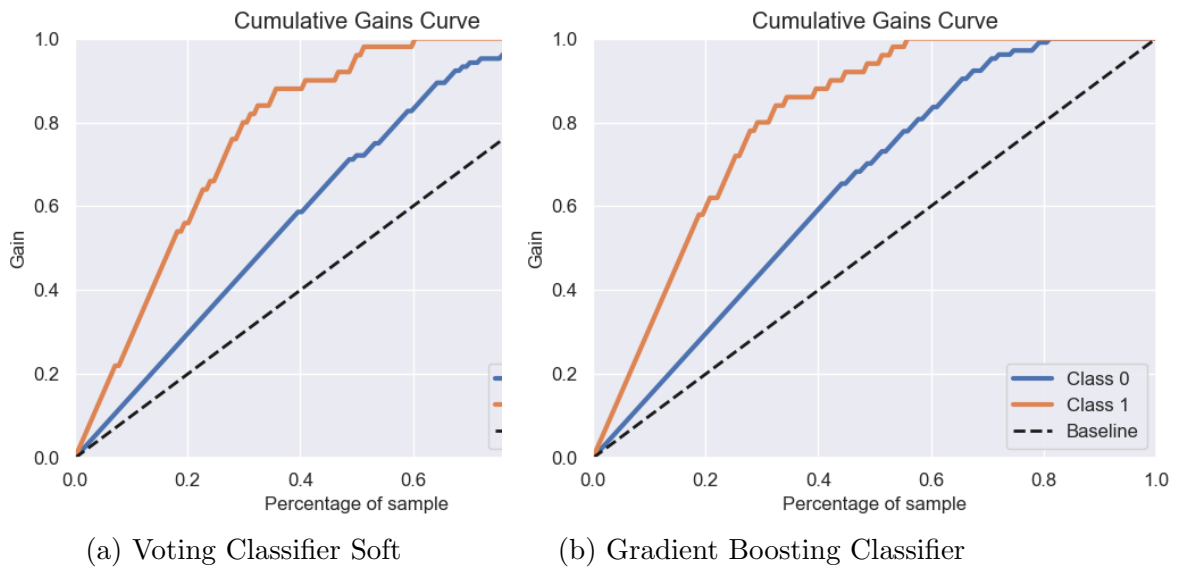


Figure 10: Cumulative Gains Curve

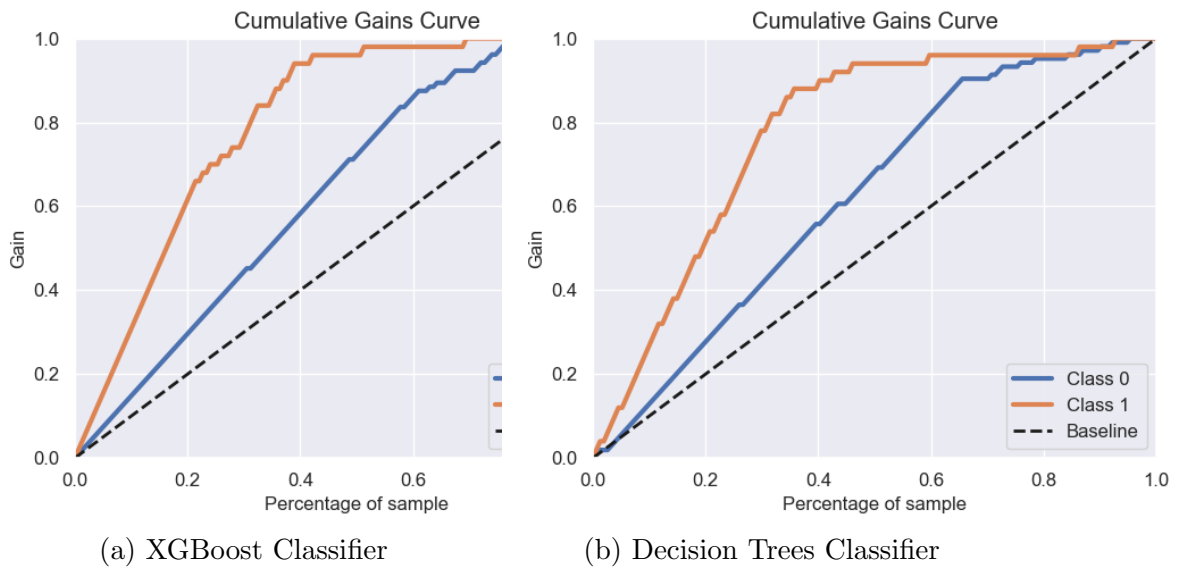


Figure 11: Cumulative Gains Curve

6 Bibliography

[1] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. The elements of statistical learning: data mining, inference, and prediction. 2009.

[2] Aurelien Geron. Hands-on machine learning with Scikit-Learn and TensorFlow: concepts, tools, and techniques to build intelligent systems. O'Reilly Media, Inc.", 2019.

[3] Lecture Notes for FYS-STK4155 at: <https://github.com/CompPhysics/MachineLearning>
access date: 16.12.2020

[4] Scikit-Learn homepage at: <https://scikit-learn.org/>

[5] The Pima Indians Diabetes Database at: <https://www.kaggle.com/uciml/pima-indians-diabetes-database>