

DATA MINING & STATISTICAL LEARNING
ISYE 7406
COURSE PROJECT

FARHAD BATMANGHELICH

11/20/2021

OUTLINE

- The dataset that is used in this project is “Diabetes Prediction” and is obtained from the following URL:
 - <https://data.world/informatics-edu/diabetes-prediction>
- This is real patient data collected from several hundred rural African American patients and is originally from biostatistics program at Vanderbilt.
- The dataset contains 390 labeled datapoints where each datapoint is a vector of 14 features describing a given patient and the label is binary “Diabetes/No Diabetes”

FEATURES

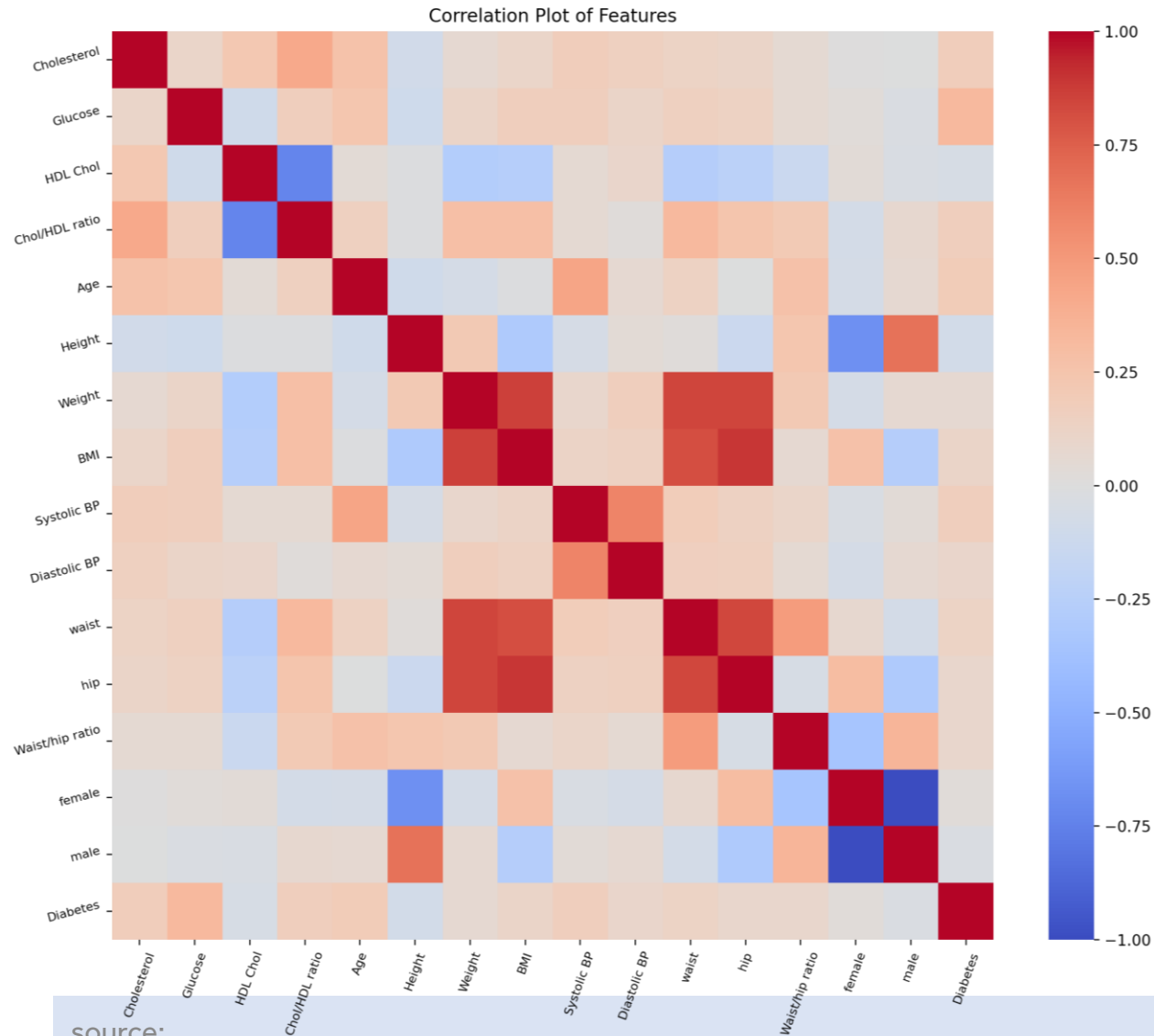
- The following table captures input features as well as the binary label

(Diabetes/No Diabetes)

Column attribute	Description
Patient number	Identifies patients by number
Cholesterol	Total cholesterol
Glucose	Fasting blood sugar
HDL	HDL or good cholesterol
Chol/HDL	Ratio of total cholesterol to good cholesterol. Desirable result is < 5
Age	All adult African Americans
Gender	162 males, 228 females
Height	In inches
Weight	In pounds (lbs)
BMI	$703 \times \text{weight (lbs)} / [\text{height(inches)}]^2$
Systolic BP	The upper number of blood pressure
Diastolic BP	The lower number of blood pressure
Waist	Measured in inches
Hip	Measured in inches
Waist/hip	Ratio is possibly a stronger risk factor for heart disease than BMI
Diabetes	Yes (60), No (330)

EXPLORATORY DATA ANALYSIS

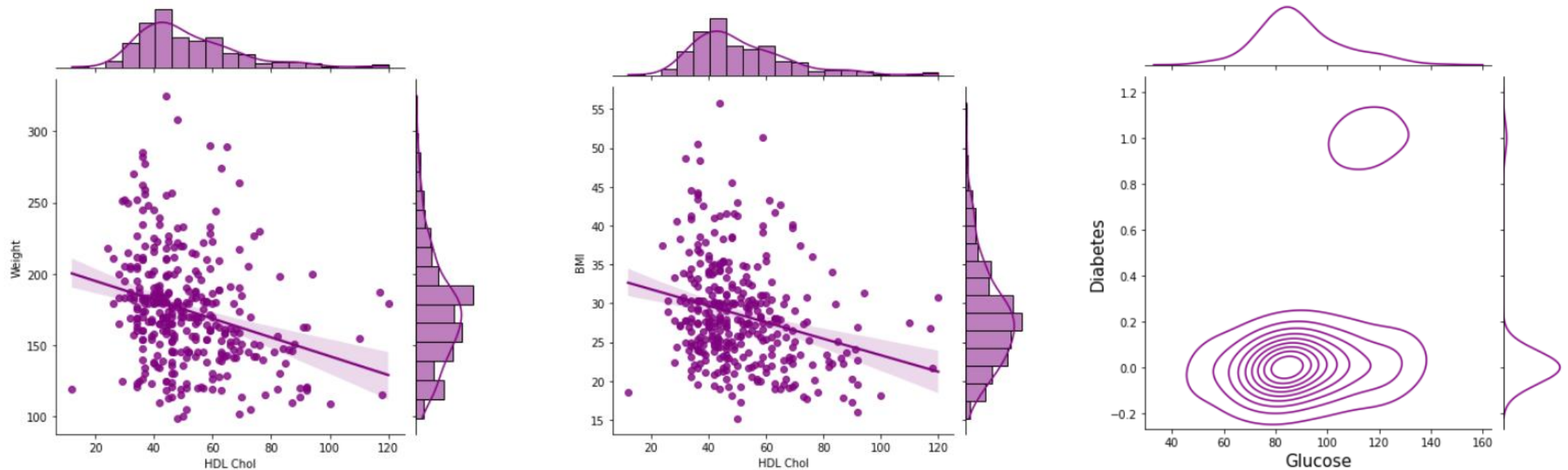
○ Correlation plot



- ❑ HDL cholesterol has a negative correlation with body mass features such as Weight, BMI, waist etc. HDL cholesterol is the so called “good” cholesterol because it helps removing other forms of cholesterol, thus higher HDL is better.
- ❑ HDL cholesterol and Height have negative correlation with having diabetes
- ❑ Glucose, cholesterol, age, Systolic BP have positive correlation with having diabetes

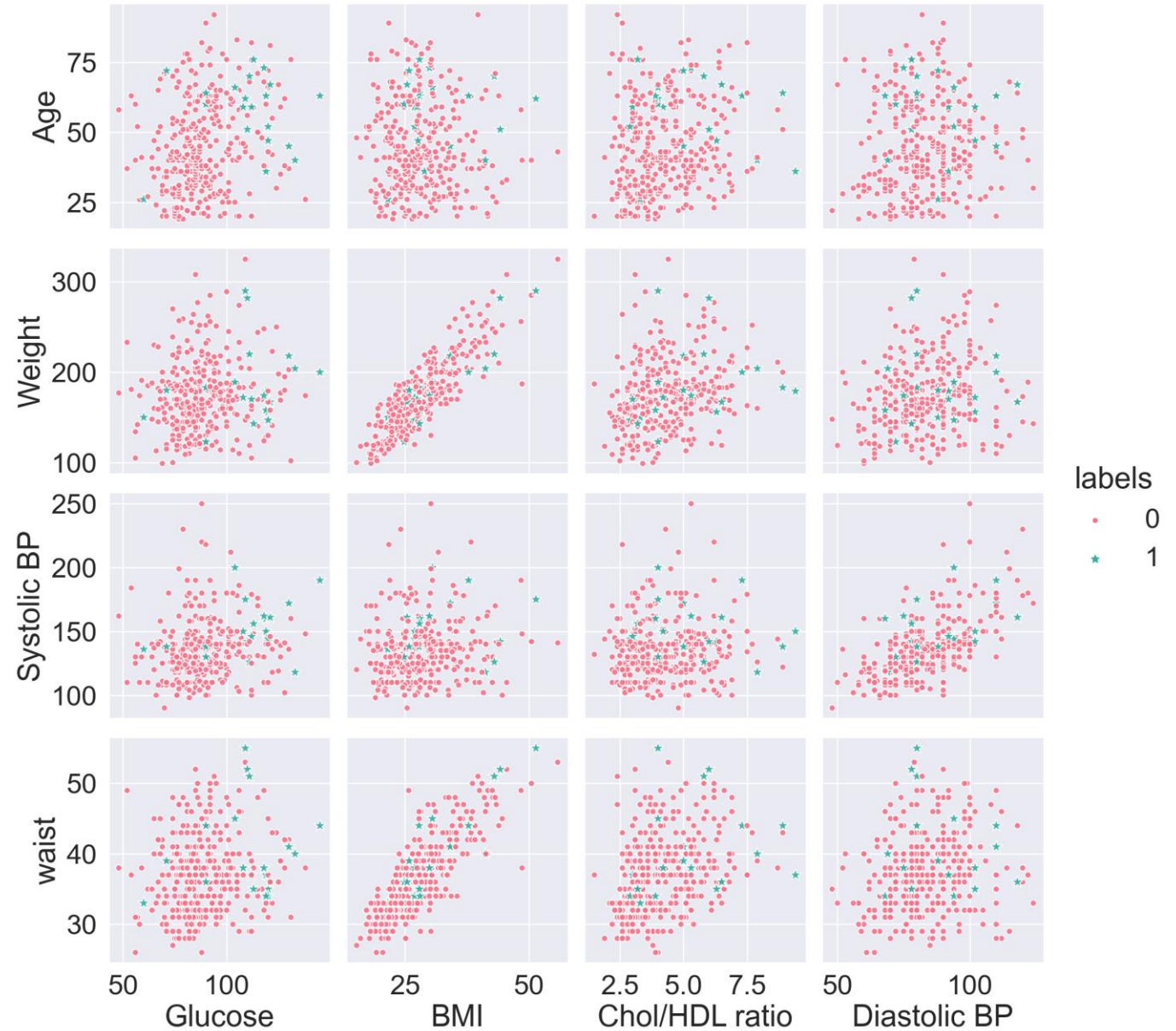
EXPLORATORY DATA ANALYSIS

- ❑ HDL cholesterol has a negative correlation (linear correlation) with body mass features such as Weight and BMI
- ❑ There are two relatively distinct distributions at low and high values of Glucose, this shows that the Glucose feature in this dataset is promising in prediction of diabetes.



MODELING

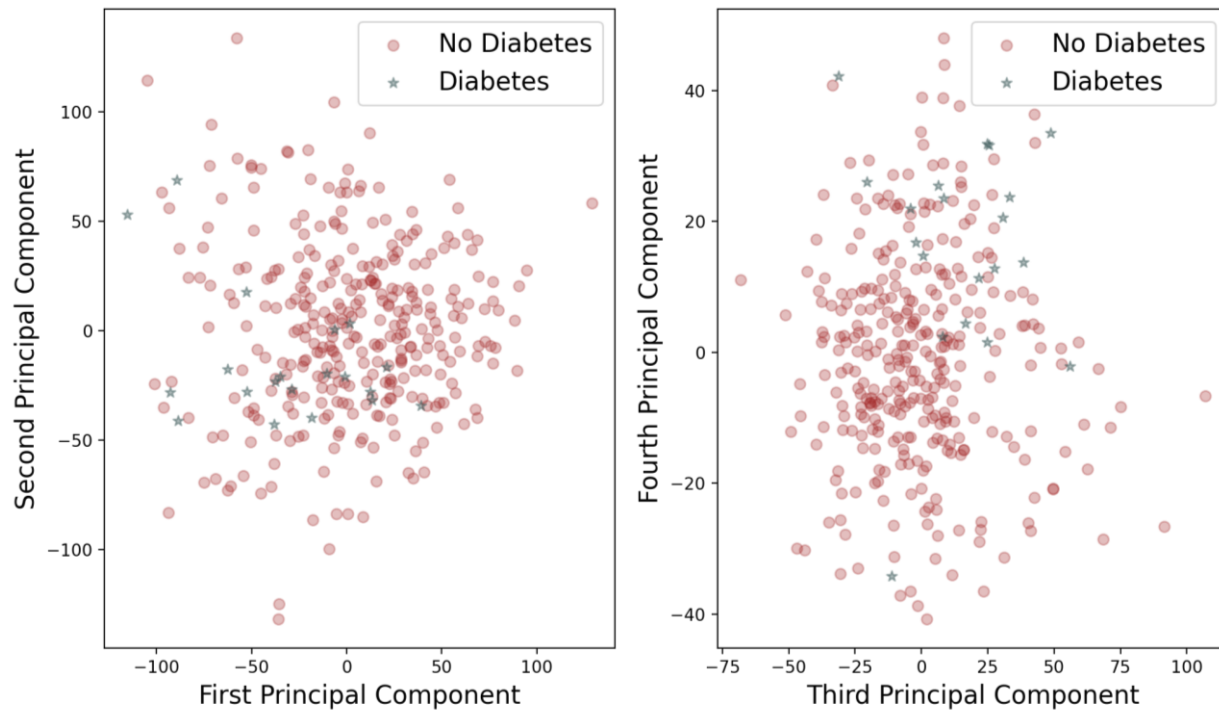
- ❑ Plotted few of the variables against each other with the label (0: no diabetes, 1: diabetes)
- ❑ Not a clear linear separation between the two classes
- ❑ Will plot PCAs to determine if there is a distinction there



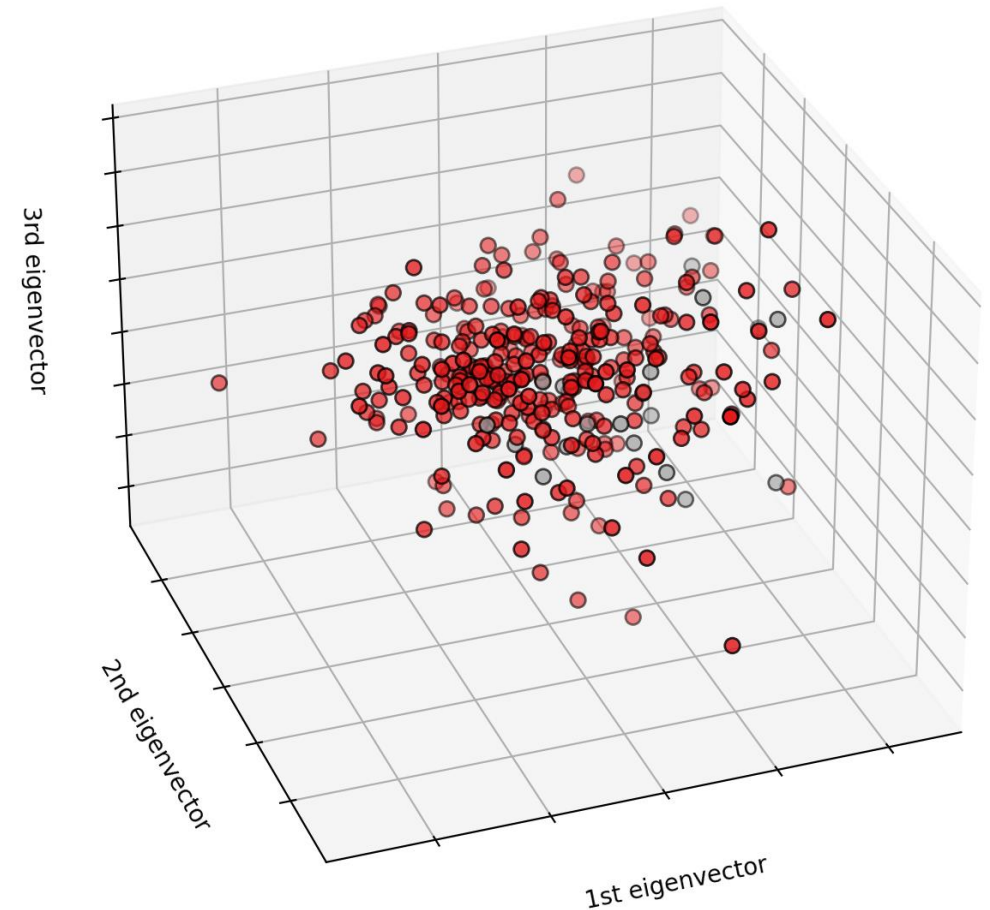
MODELING

- ❑ Seems that in the principal component space, there is still no linear boundary between classes

PCA 2D



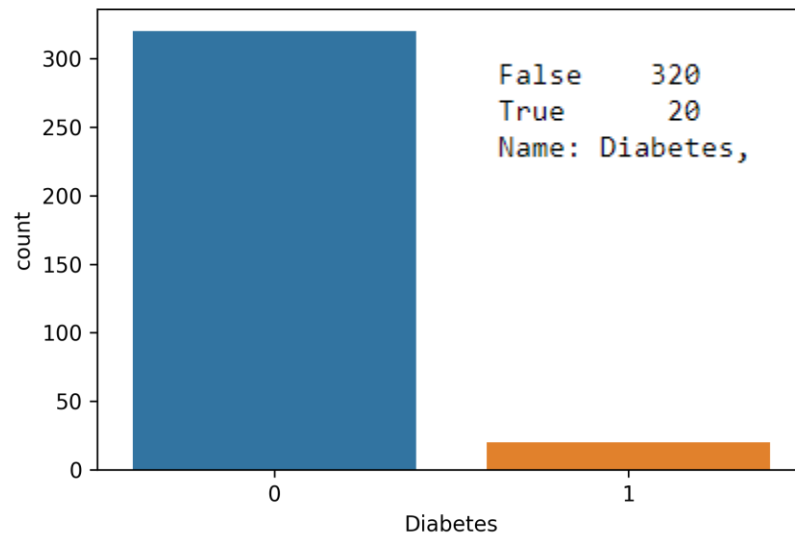
PCA 3D



MODELING

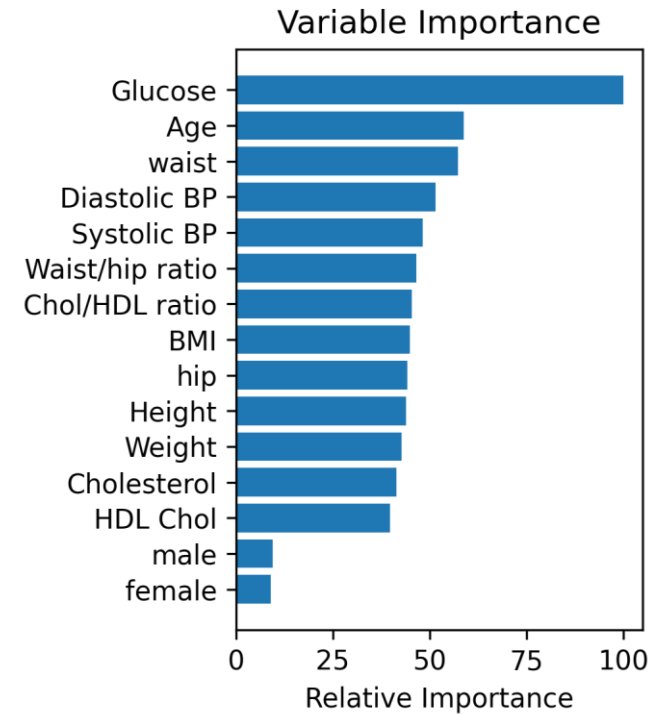
❑ As can be seen from the previous results:

- This data is highly imbalanced



- The two classes are not linearly separable

❑ Seems that glucose, age and waist are important in predicting diabetes

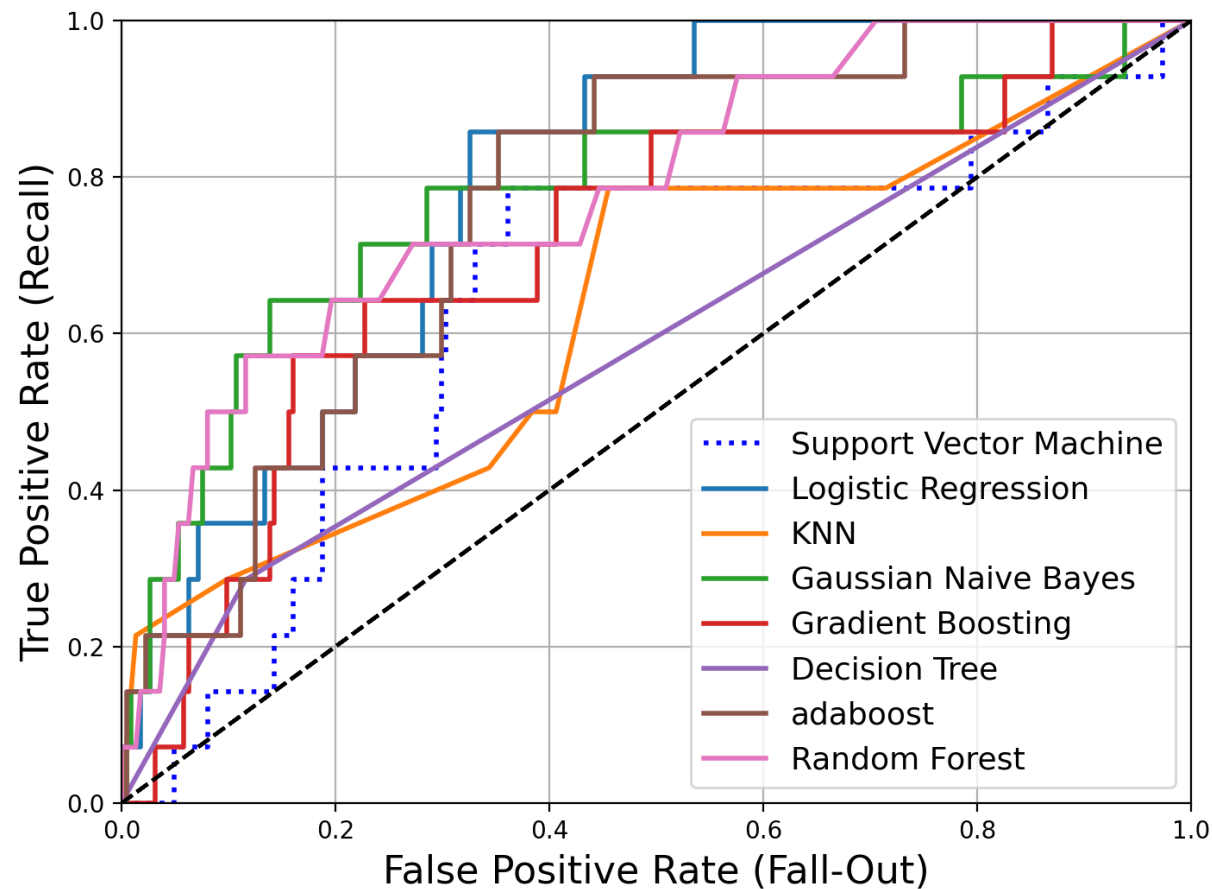


MODELING-dealing with imbalance in classes

- ❑ I tried “imblearn” library to deal with scarcity of class “1” i.e. Diabetes
- ❑ These methods are based on:
 - Under-sampling
 - Over-sampling
 - Combination of the above
- ❑ I used the **RepeatedEditedNearestNeighbours** class which is an under-sampling method on majority class i.e. No Diabetes:
 - Under-sample based on the repeated edited nearest neighbor method

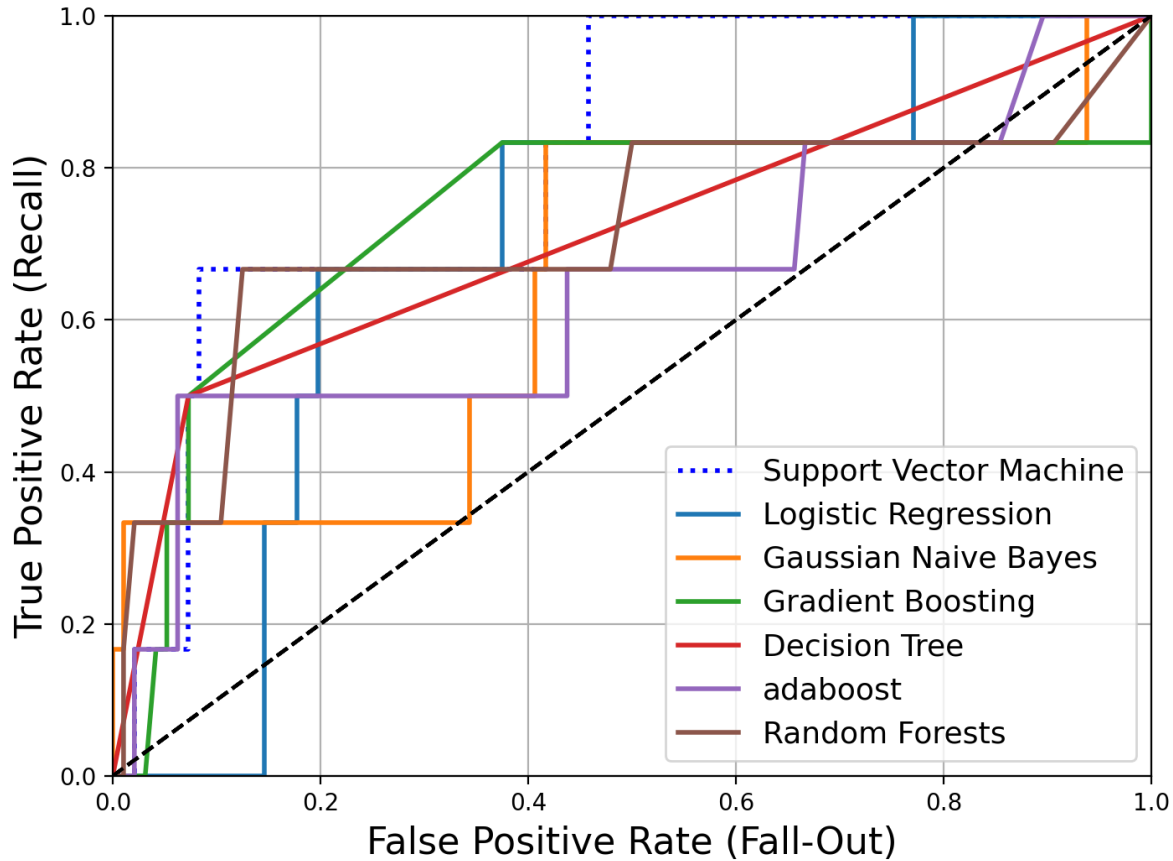
https://imbalanced-learn.org/stable/references/generated/imblearn.under_sampling.RepeatedEditedNearestNeighbours.html

MODEL PERFORMANCE ON TRAINING DATA



Model	Classification Report															
Support Vector Machine	<div>Confusion Matrix:</div> <div>[[224 0] [14 0]]</div> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>0.94</td><td>1.00</td><td>0.97</td><td>224</td></tr><tr><td>1</td><td>0.00</td><td>0.00</td><td>0.00</td><td>14</td></tr></table>		precision	recall	f1-score	support	0	0.94	1.00	0.97	224	1	0.00	0.00	0.00	14
	precision	recall	f1-score	support												
0	0.94	1.00	0.97	224												
1	0.00	0.00	0.00	14												
Logistic Regression	<div>Confusion Matrix:</div> <div>[[179 45] [2 12]]</div> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>0.99</td><td>0.80</td><td>0.88</td><td>224</td></tr><tr><td>1</td><td>0.21</td><td>0.86</td><td>0.34</td><td>14</td></tr></table>		precision	recall	f1-score	support	0	0.99	0.80	0.88	224	1	0.21	0.86	0.34	14
	precision	recall	f1-score	support												
0	0.99	0.80	0.88	224												
1	0.21	0.86	0.34	14												
KNN	<div>Confusion Matrix:</div> <div>[[202 22] [5 9]]</div> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>0.98</td><td>0.90</td><td>0.94</td><td>224</td></tr><tr><td>1</td><td>0.29</td><td>0.64</td><td>0.40</td><td>14</td></tr></table>		precision	recall	f1-score	support	0	0.98	0.90	0.94	224	1	0.29	0.64	0.40	14
	precision	recall	f1-score	support												
0	0.98	0.90	0.94	224												
1	0.29	0.64	0.40	14												
Gaussian Naïve Bayes	<div>Confusion Matrix:</div> <div>[[211 13] [6 8]]</div> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>0.97</td><td>0.94</td><td>0.96</td><td>224</td></tr><tr><td>1</td><td>0.38</td><td>0.57</td><td>0.46</td><td>14</td></tr></table>		precision	recall	f1-score	support	0	0.97	0.94	0.96	224	1	0.38	0.57	0.46	14
	precision	recall	f1-score	support												
0	0.97	0.94	0.96	224												
1	0.38	0.57	0.46	14												
Gradient Boosting	<div>Confusion Matrix:</div> <div>[[223 1] [0 14]]</div> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>1.00</td><td>1.00</td><td>1.00</td><td>224</td></tr><tr><td>1</td><td>0.93</td><td>1.00</td><td>0.97</td><td>14</td></tr></table>		precision	recall	f1-score	support	0	1.00	1.00	1.00	224	1	0.93	1.00	0.97	14
	precision	recall	f1-score	support												
0	1.00	1.00	1.00	224												
1	0.93	1.00	0.97	14												
Decision Tree	<div>Confusion Matrix:</div> <div>[[215 9] [0 14]]</div> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>1.00</td><td>0.96</td><td>0.98</td><td>224</td></tr><tr><td>1</td><td>0.61</td><td>1.00</td><td>0.76</td><td>14</td></tr></table>		precision	recall	f1-score	support	0	1.00	0.96	0.98	224	1	0.61	1.00	0.76	14
	precision	recall	f1-score	support												
0	1.00	0.96	0.98	224												
1	0.61	1.00	0.76	14												
adaboost	<div>Confusion Matrix:</div> <div>[[205 19] [0 14]]</div> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>1.00</td><td>0.92</td><td>0.96</td><td>224</td></tr><tr><td>1</td><td>0.42</td><td>1.00</td><td>0.60</td><td>14</td></tr></table>		precision	recall	f1-score	support	0	1.00	0.92	0.96	224	1	0.42	1.00	0.60	14
	precision	recall	f1-score	support												
0	1.00	0.92	0.96	224												
1	0.42	1.00	0.60	14												
Random Forest	<div>Confusion Matrix:</div> <div>[[204 20] [0 14]]</div> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>1.00</td><td>0.91</td><td>0.95</td><td>224</td></tr><tr><td>1</td><td>0.41</td><td>1.00</td><td>0.58</td><td>14</td></tr></table>		precision	recall	f1-score	support	0	1.00	0.91	0.95	224	1	0.41	1.00	0.58	14
	precision	recall	f1-score	support												
0	1.00	0.91	0.95	224												
1	0.41	1.00	0.58	14												

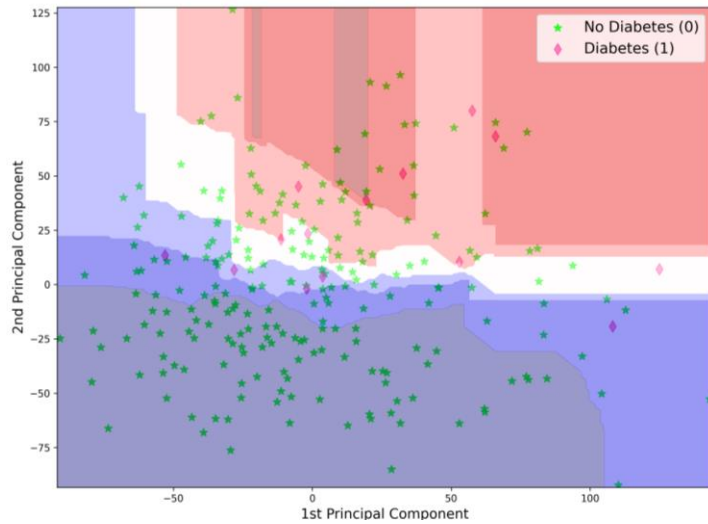
MODEL PERFORMANCE ON TEST DATA



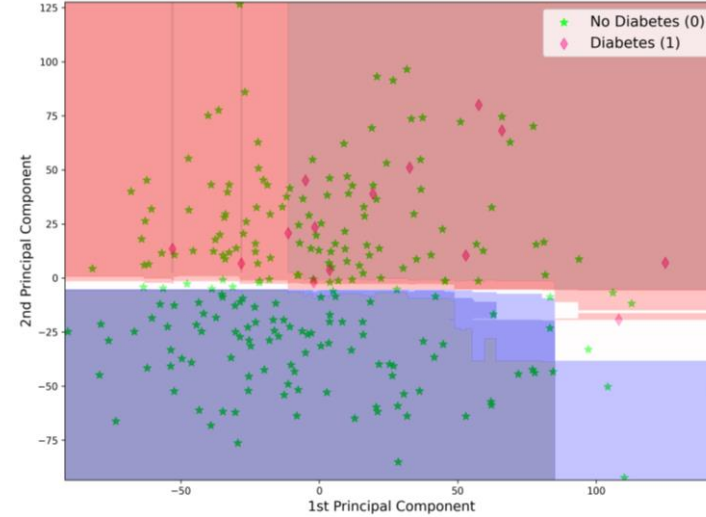
Model	Classification Report
Support Vector Machine	<div>Confusion Matrix:</div> <div><div><div>[[96 0] [6 0]]</div><div>precision</div><div>recall</div><div>f1-score</div><div>support</div></div><div><div>00.94</div><div>1.00</div><div>0.97</div><div>96</div></div><div><div>10.00</div><div>0.00</div><div>0.00</div><div>6</div></div></div>
Logistic Regression	<div>Confusion Matrix:</div> <div><div><div>[[73 23] [1 5]]</div><div>precision</div><div>recall</div><div>f1-score</div><div>support</div></div><div><div>00.99</div><div>0.76</div><div>0.86</div><div>96</div></div><div><div>10.18</div><div>0.83</div><div>0.29</div><div>6</div></div></div>
KNN	<div>Confusion Matrix:</div> <div><div><div>[[85 11] [2 4]]</div><div>precision</div><div>recall</div><div>f1-score</div><div>support</div></div><div><div>00.98</div><div>0.89</div><div>0.93</div><div>96</div></div><div><div>10.27</div><div>0.67</div><div>0.38</div><div>6</div></div></div>
Gaussian Naïve Bayes	<div>Confusion Matrix:</div> <div><div><div>[[89 7] [3 3]]</div><div>precision</div><div>recall</div><div>f1-score</div><div>support</div></div><div><div>00.97</div><div>0.93</div><div>0.95</div><div>96</div></div><div><div>10.30</div><div>0.50</div><div>0.37</div><div>6</div></div></div>
Gradient Boosting	<div>Confusion Matrix:</div> <div><div><div>[[95 1] [5 1]]</div><div>precision</div><div>recall</div><div>f1-score</div><div>support</div></div><div><div>00.95</div><div>0.99</div><div>0.97</div><div>96</div></div><div><div>10.50</div><div>0.17</div><div>0.25</div><div>6</div></div></div>
Decision Tree	<div>Confusion Matrix:</div> <div><div><div>[[91 5] [2 4]]</div><div>precision</div><div>recall</div><div>f1-score</div><div>support</div></div><div><div>00.98</div><div>0.95</div><div>0.96</div><div>96</div></div><div><div>10.44</div><div>0.67</div><div>0.53</div><div>6</div></div></div>
adaboost	<div>Confusion Matrix:</div> <div><div><div>[[86 10] [3 3]]</div><div>precision</div><div>recall</div><div>f1-score</div><div>support</div></div><div><div>00.97</div><div>0.90</div><div>0.93</div><div>96</div></div><div><div>10.23</div><div>0.50</div><div>0.32</div><div>6</div></div></div>
Random Forest	<div>Confusion Matrix:</div> <div><div><div>[[89 7] [2 4]]</div><div>precision</div><div>recall</div><div>f1-score</div><div>support</div></div><div><div>00.98</div><div>0.93</div><div>0.95</div><div>96</div></div><div><div>10.36</div><div>0.67</div><div>0.47</div><div>6</div></div></div>

DECISION BOUNDARIES LEARNED ON PCA SPACE FOR TRAINING

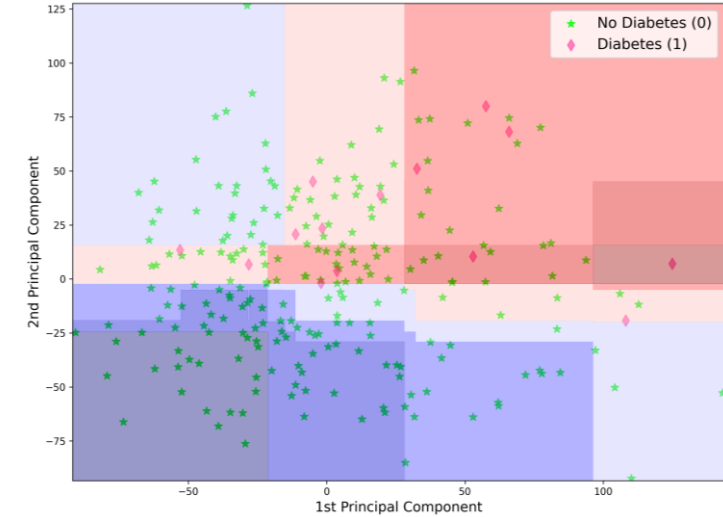
KNN



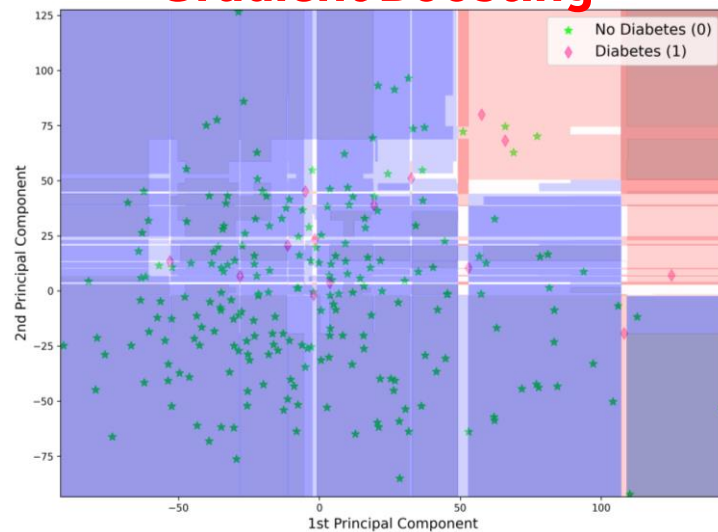
Random forests



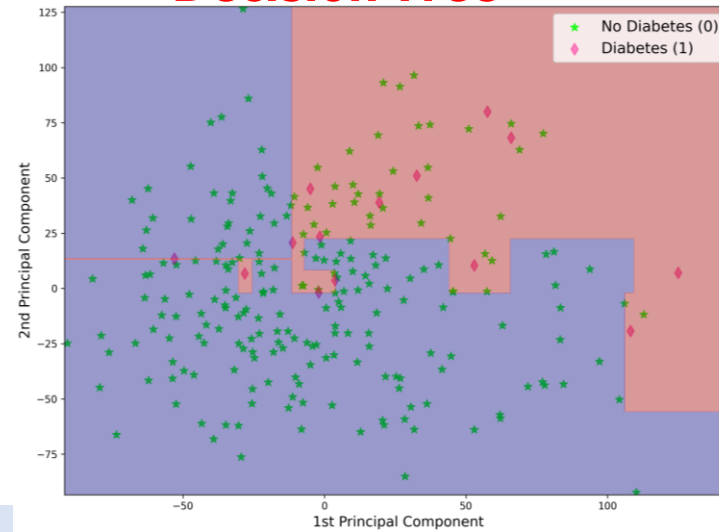
Adaboost



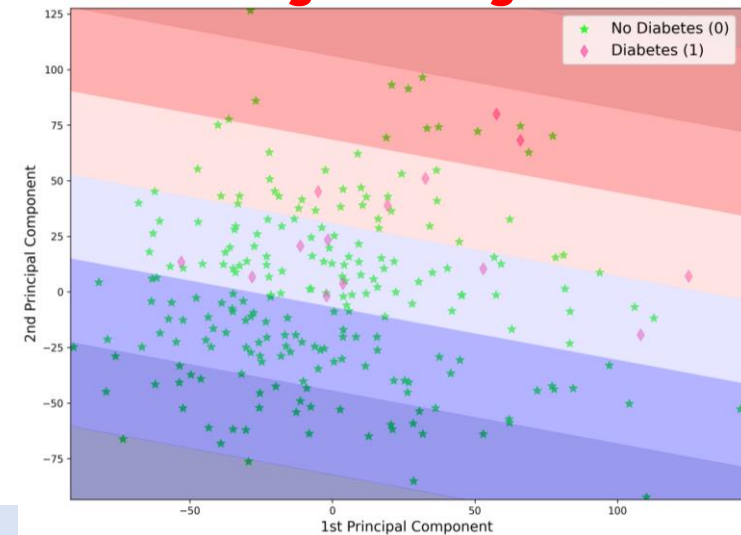
Gradient Boosting



Decision Tree



Logistic Regression



CONCLUSIONS

- ❑ The imbalance in distribution of positive class is a major issue in fitting a classification model to diagnostics datasets such as diabetes, cancer etc.
- ❑ Model performance can be boosted with under-sampling, over-sampling or both on either majority or minority class
- ❑ Classes seemed to not be linearly separable
- ❑ Ensemble models such as random forests and adaboost seemed to be able to distinguish between classes given their precision, recall and f1-scores.
- ❑ Imbalance in dataset should be address if one of these models were to be deployed to be used by physicians.