

# Lab Assignment 5

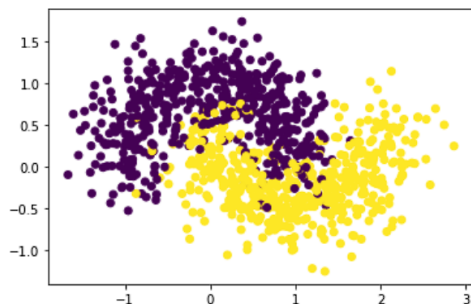
## Prakhar Gupta

### B21AI027

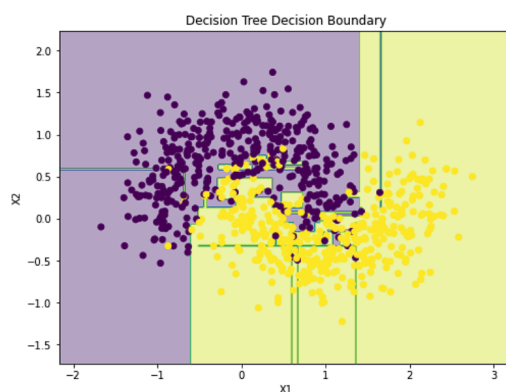
#### Question 1:

##### Part 1-

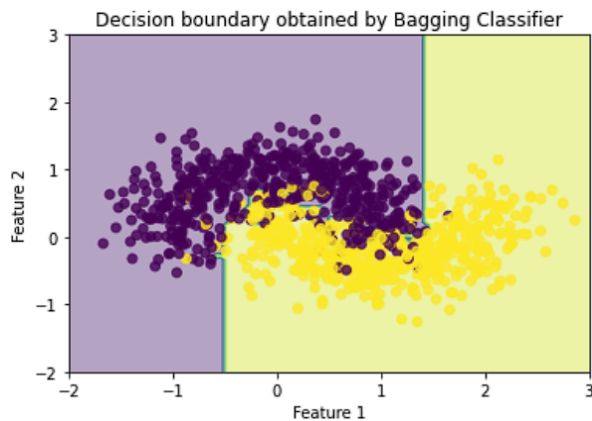
- Made **X,y** using **sklearn make\_moons** function(`random_state=42`, `noise=0.3`, `n_samples=1000`)
- Converted **X,y** to **df** and then used **df.describe()**
- Checked **df.isnull().sum()**
- Plotted X,y



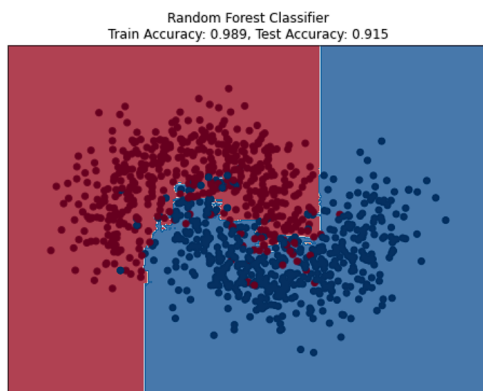
- Used `train_test_split` to split train:test in **80:20**
- Fitted Decision Tree Classifier



- Plotted DT boundary
- Fitted best params for DT using **GridSearchCV**
- Trained a Bagging Classifier
- Plotted Decision boundary



- Trained a RandomForestClassifier
- Plotted Decision Boundary



- Best accuracy order of Train\_data - **Random Forest > Bagging Classifier > Decision Tree**
- Best accuracy order of Test\_data - **Random Forest > Bagging Classifier > Decision Tree**
- We can see Decision Tree has **high bias**, Bagging Tree is somewhat **overfitting** and **Random Forest is the best**
- Varied the number of estimators for the BaggingClassifier and RandomForestClassifier and plotted Decision Boundaries
- From The Accuracy we can clearly see that it sometimes up and down but on long run it clearly increases
- The best accuracy we get are:
- bst\_train\_randforest\_accuracy: 0.99375
- bst\_test\_randforest\_accuracy: 0.915
- bst\_train\_bag\_tree\_accuracy: 0.99375
- bst\_test\_bag\_tree\_accuracy: 0.91
- We can see the randomforest and bag\_tree accuracy get **equal at high estimators** value in train set
- We can see the randomforest accuracy is greater than bag\_tree accuracy at high estimators value in test set

- From plots of both bag\_tree and random forest we can see that the boundary becomes more and more **complex** as n\_estimators value increases

## Part 2-

- Implemented Bagging algorithm from scratch
- Applied scratch bagging algorithm with **n\_estimators = 10**
- Trained the Bagging algo got its accuracy `0.92`
- Average accuracy on trees is `0.9`
- The average accuracy of all the trees is 0.9 whereas for bagging it is 0.92
- Bagging accuracy > Average accuracy of all trees

## Question 2:

### Part 1-

- Trained AdaBoostClassifier with n\_estimators=100

### Part 2-

- Trained a XGBClassifier with n\_estimators=100, and subsample=0.7

### Part 3-

- Printed accuracy of both AdaBoostClassifier and XGBClassifier on train and test data

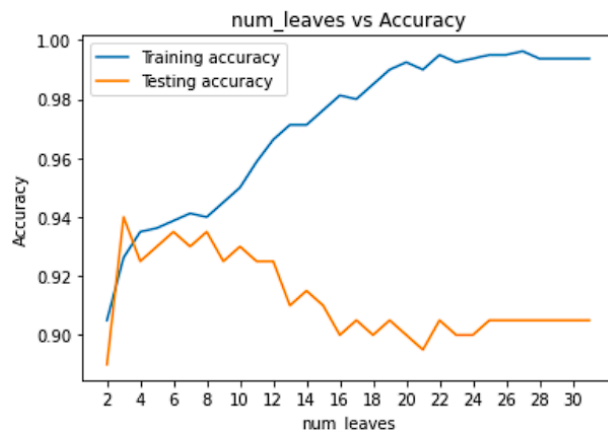
```
AdaBoost Classifier Accuracy on training set: 0.9475
AdaBoost Classifier Accuracy on test set: 0.91
XGBoost Classifier Accuracy on training set: 0.9925
XGBoost Classifier Accuracy on test set: 0.905
```

### Part 4-

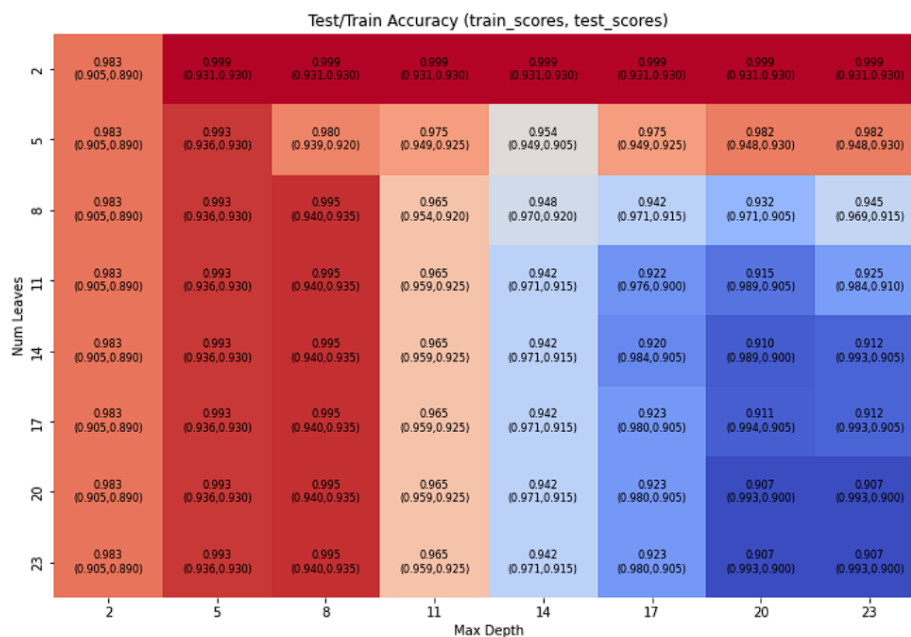
- Trained a LightGBM model
- Varied different values for num\_leaves

### Part 5-

- Analysed relation between max\_depth when num\_leaves is given to the LightGBM model
- Plotted accuracy for various num\_leaves value



- Plotted heatmap where axis are max\_depth and num\_leaves and the y value is test/train accuracy



- From num\_leaves vs Accuracy plot we can see the model starts overfitting after **num\_leaves=8** and **max\_depth=8**
- From heatmap also we can see the same result

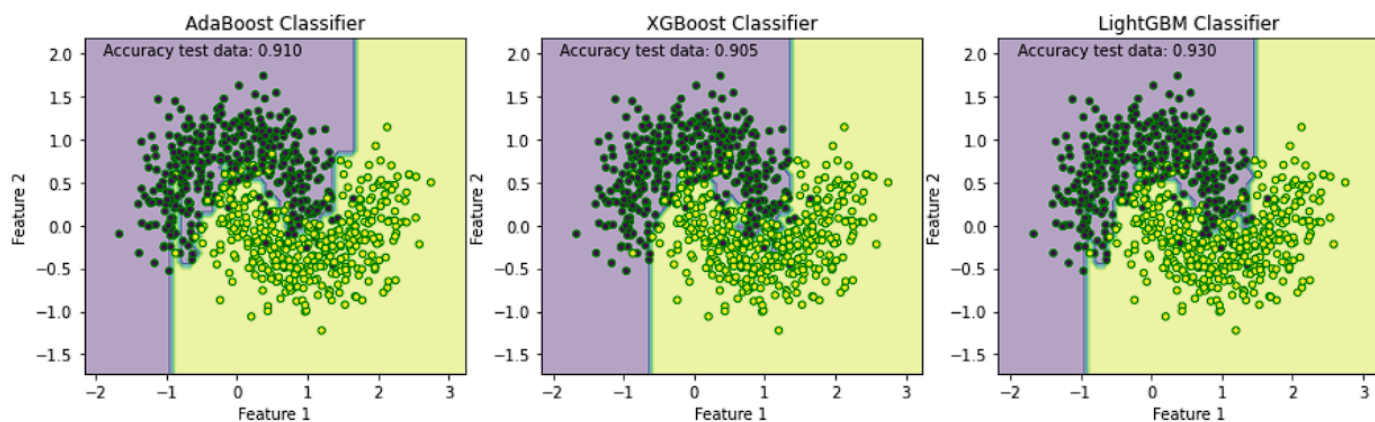
## Part 6-

- num\_leaves**: This decides how many leaves will be in a tree. If we increase num\_leaves, it can make the prediction better.
- max\_depth**: High depth can also lead to overfitting. But it will increase the accuracy
- min\_data\_in\_leaf**: This parameter sets the minimum number of samples required to be at a leaf node. If we increase this value, we will have less overfitting for each leaf, which can help to prevent overfitting.

- **feature\_fraction**: This decides how many features we want to use in each tree. If we lower this value, we will have less overfitting, which can help us to prevent overfitting.
- **lambda\_l1 and lambda\_l2**: These parameters decide L1 and L2 regularization, respectively. If we set non-zero values for lambda\_l1 and lambda\_l2, it can help us to prevent overfitting.

## Part 7-

- Plot the decision boundaries for all the 3 models and compared their performance



- AdaBoost Classifier Accuracy on training set: 0.9475 and on test set: 0.91
- XGBoost Classifier Accuracy on training set: 0.9925 and on test set: 0.905
- LightGBM Classifier Accuracy on training set: 0.9375 and on test set: 0.93
- We can clearly see that accuracy/performance is **LightGBM>AdaBoost>XGBoost**

## Question 3:

### Part 1-

- Train a Bayes classification model and also tuned the hyperparameter(i.e. var\_smoothing)

```
Best parameters found: {'var_smoothing': 1e-19}
```

```
Gaussian Naive Bayes Classifier Accuracy on training set: 0.86
```

```
Gaussian Naive Bayes Classifier Accuracy on test set: 0.82
```

- The best models till now are RandomForestClassifier, LightGBM Classifier, AdaBoost Classifier

- Made a VotingClassifier from sklearn using (Bayes Classification model, RandomForestClassifier, LightGBM Classifier, AdaBoost Classifier)
- Trained VotingClassifier
- `Voting Classifier Accuracy on training set: 0.95125`  
`Voting Classifier Accuracy on test set: 0.915`
- Voting Classifier Accuracy on training set: 0.9475 and on test set: 0.91
- Average Accuracy on training set by AdaBoost, Randomforst, LightGBM: 0.95833 and and on testing set: 0.91833
- AdaBoost Classifier Accuracy on training set: 0.9475 and on test set: 0.91
- Random Forest Accuracy on training set: 0.95125 and on test set: 0.915
- LightGBM Classifier Accuracy on training set: 0.9375 and on test set: 0.93
- In comparison to the three model Voting Classifier has both train\_accuracy and test\_accuracy good