



**Department of Computer Science and Engineering (Data Science)**

**Subject: Machine Learning – I (DJ19DSC402)**

**AY: 2022-23**

**Jhanvi Parekh**

**60009210033**

**CSE(Data Science)**

**Experiment 7**

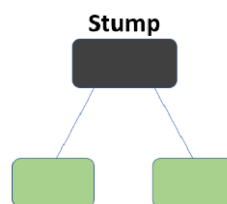
**(AdaBoost)**

**Aim:** Evaluate the performance of boosting algorithm (AdaBoost) with different base learners and hyperparameter tuning.

**Theory:**

Boosting algorithms improve the prediction power by **converting a number of weak learners to strong learners**. The principle behind boosting algorithms is first built a model on the training dataset, then a second model is built to rectify the errors present in the first model. This procedure is continued until and unless the errors are minimized, and the dataset is predicted correctly. Let's take an example to understand this, suppose you built a decision tree algorithm on the Titanic dataset and from there you get an accuracy of 80%. After this, you apply a different algorithm and check the accuracy and it comes out to be 75% for KNN and 70% for Linear Regression. The accuracy differs when we built a different model on the same dataset. But what if we use combinations of all these algorithms for making the final prediction? We'll get more accurate results by taking the average of results from these models. We can increase the prediction power in this way.

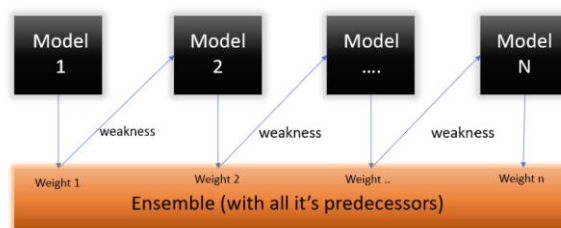
AdaBoost also called Adaptive Boosting is a technique in Machine Learning used as an Ensemble Method. The most common algorithm used with AdaBoost is decision trees with one level that means with Decision trees with only 1 split. These trees are also called **Decision Stumps**.





## Department of Computer Science and Engineering (Data Science)

Algorithm builds a model and gives equal weights to all the data points. It then assigns higher weights to points that are wrongly classified. Now all the points which have higher weights are given more importance in the next model. It will keep training models until and unless a lower error is received.



**Step 1** – The Image is shown below is the actual representation of our dataset. Since the target column is binary it is a classification problem. First of all these data points will be assigned some weights. Initially, all the weights will be equal.

Row No.	Gender	Age	Income	Illness	Sample Weights
1	Male	41	40000	Yes	1/5
2	Male	54	30000	No	1/5
3	Female	42	25000	No	1/5
4	Female	40	60000	Yes	1/5
5	Male	46	50000	Yes	1/5

The formula to calculate the sample weights is:

$$w(x_i, y_i) = \frac{1}{N}, i = 1, 2, \dots, n$$

Where N is the total number of datapoints. since we have 5 data points so the sample weights assigned will be 1/5.

**Step 2** – We start by seeing how well “Gender” classifies the samples and will see how the variables (Age, Income) classifies the samples.

We’ll create a decision stump for each of the features and then calculate the **Gini Index** of each tree. The tree with the lowest Gini Index will be our first stump.

Here in our dataset let’s say **Gender** has the lowest gini index so it will be our first stump.

**Step 3** – Calculate the “Amount of Say” or “Importance” or “Influence” for this classifier in classifying the datapoints using this formula:

$$\frac{1}{2} \log \frac{1 - \text{Total Error}}{\text{Total Error}}$$

The total error is nothing, but the summation of all the sample weights of misclassified data points.

Here in our dataset let’s assume there is 1 wrong output, so our total error will be 1/5, and alpha(performance of the stump) will be:



## Department of Computer Science and Engineering (Data Science)

$$\text{Performance of the stump} = \frac{1}{2} \log_e \left( \frac{1 - \text{Total Error}}{\text{Total Error}} \right)$$

$$\alpha = \frac{1}{2} \log_e \left( \frac{1 - \frac{1}{5}}{\frac{1}{5}} \right)$$

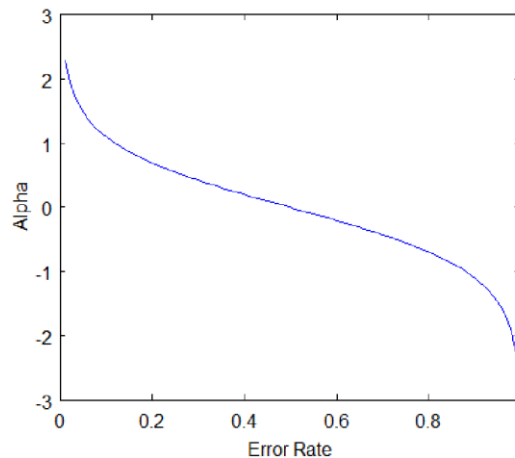
$$\alpha = \frac{1}{2} \log_e \left( \frac{0.8}{0.2} \right)$$

$$\alpha = \frac{1}{2} \log_e(4) = \frac{1}{2} * (1.38)$$

$$\alpha = 0.69$$

**Note:** Total error will always be between 0 and 1.

0 Indicates perfect stump and 1 indicates horrible stump.



From the graph above we can see that when there is no misclassification then we have no error (Total Error = 0), so the “amount of say (alpha)” will be a large number.

When the classifier predicts half right and half wrong then the Total Error = 0.5 and the importance (amount of say) of the classifier will be 0. If all the samples have been incorrectly classified then the error will be very high (approx. to 1) and hence our alpha value will be a negative integer.

**Step 4** –We need to update the weights because if the same weights are applied to the next model, then the output received will be the same as what was received in the first model. The wrong predictions will be given more weight whereas the correct predictions weights will be decreased. Now when we build our next model after updating the weights, more preference will be given to the points with higher weights. After finding the importance of the classifier and total error we need to finally update the weights and for this, we use the following formula:

$$\text{New sample weight} = \text{old weight} * e^{\pm \text{Amount of say } (\alpha)}$$

The amount of say (alpha) will be **negative** when the sample is **correctly classified**.

The amount of say (alpha) will be **positive** when the sample is **miss-classified**.

There are four correctly classified samples and 1 wrong, here the **sample weight** of that datapoint is 1/5 and the **amount of say/performance of the stump** of **Gender** is 0.69.



## Department of Computer Science and Engineering (Data Science)

$$\text{New sample weight} = \frac{1}{5} * \exp(-0.69)$$

$$\text{New sample weight} = 0.2 * 0.502 = 0.1004$$

For *wrongly classified* samples the updated weights will be:

$$\text{New sample weight} = \frac{1}{5} * \exp(0.69)$$

$$\text{New sample weight} = 0.2 * 1.994 = 0.3988$$

**Note:** See the sign of alpha when I am putting the values, the **alpha is negative** when the data point is correctly classified, and this *decreases the sample weight* from 0.2 to 0.1004. It is **positive** when there is **misclassification**, and this will *increase the sample weight* from 0.2 to 0.3988

Row No.	Gender	Age	Income	Illness	Sample Weights	New Sample Weights
1	Male	41	40000	Yes	1/5	0.1004
2	Male	54	30000	No	1/5	0.1004
3	Female	42	25000	No	1/5	0.1004
4	Female	40	60000	Yes	1/5	0.3988
5	Male	46	50000	Yes	1/5	0.1004

We know that the total sum of the sample weights must be equal to 1 but here if we sum up all the new sample weights, we will get 0.8004. To bring this sum equal to 1 we will normalize these weights by dividing all the weights by the total sum of updated weights that is 0.8004. So, after normalizing the sample weights we get this dataset and now the sum is equal to 1.

Row No.	Gender	Age	Income	Illness	Sample Weights	New Sample Weights
1	Male	41	40000	Yes	1/5	0.1004/0.8004 =0.1254
2	Male	54	30000	No	1/5	0.1004/0.8004 =0.1254
3	Female	42	25000	No	1/5	0.1004/0.8004 =0.1254
4	Female	40	60000	Yes	1/5	0.3988/0.8004 =0.4982
5	Male	46	50000	Yes	1/5	0.1004/0.8004 =0.1254

**Step 5** – Now we need to make a new dataset to see if the errors decreased or not. For this we will remove the “sample weights” and “new sample weights” column and then based on the “new sample weights” we will divide our data points into buckets.

Row No.	Gender	Age	Income	Illness	New Sample Weights	Buckets
1	Male	41	40000	Yes	0.1004/0.8004= 0.1254	0 to 0.1254
2	Male	54	30000	No	0.1004/0.8004= 0.1254	0.1254 to 0.2508
3	Female	42	25000	No	0.1004/0.8004= 0.1254	0.2508 to 0.3762
4	Female	40	60000	Yes	0.3988/0.8004= 0.4982	0.3762 to 0.8744
5	Male	46	50000	Yes	0.1004/0.8004= 0.1254	0.8744 to 0.9998

**Step 6** – We are almost done, now what the algorithm does is selects random numbers from 0-1. Since incorrectly classified records have higher sample weights, the probability to select those records is very high. Suppose the 5 random numbers our algorithm take is 0.38,0.26,0.98,0.40,0.55.

Now we will see where these random numbers fall in the bucket and according to it, we'll make our new dataset shown below.



### Department of Computer Science and Engineering (Data Science)

Row No.	Gender	Age	Income	Illness
1	Female	40	60000	Yes
2	Male	54	30000	No
3	Female	42	25000	No
4	Female	40	60000	Yes
5	Female	40	60000	Yes

This comes out to be our new dataset and we see the datapoint which was wrongly classified has been selected 3 times because it has a higher weight.

**Step 9** – Now this act as our new dataset and we need to repeat all the above steps i.e.

1. Assign **equal weights** to all the datapoints
2. Find the stump that does the **best job classifying** the new collection of samples by finding their Gini Index and selecting the one with the lowest Gini index
3. Calculate the **"Amount of Say"** and **"Total error"** to update the previous sample weights.
4. Normalize the new sample weights.

Iterate through these steps until and unless a low training error is achieved.

Suppose with respect to our dataset we have constructed 3 decision trees (DT1, DT2, DT3) in a **sequential manner**. If we send our **test data** now it will pass through all the decision trees and finally, we will see which class has the majority, and based on that we will do predictions for our test dataset.

### Lab Assignments to complete in this session:

Use the given dataset and perform the following tasks:

#### Dataset 1: Synthetic dataset

**Dataset 2: CreditcardFraud.csv:** The dataset contains transactions made by credit cards in September 2013 by European cardholders. This dataset presents transactions that occurred in two days, where we have 492 frauds out of 284,807 transactions. The dataset is highly unbalanced, the positive class (frauds) account for 0.172% of all transactions. It contains only numerical input variables which are the result of a PCA transformation. Unfortunately, due to confidentiality issues, the original features and more background information about the data are not provided. Features V1, V2, ... V28 are the principal components obtained with PCA, the only features which have not been transformed with PCA are 'Time' and 'Amount'. Feature 'Time' contains the seconds elapsed between each transaction and the first transaction in the dataset. The feature 'Amount' is the transaction Amount, this feature can be used for example-dependant cost-sensitive learning. Feature 'Class' is the response variable and it takes value 1 in case of fraud and 0 otherwise.

1. Implement Decision Tree classifier and Logistic Regression on Dataset 1 using K fold cross validation and compare the results with AdaBoost classifier with base learner as Decision tree and Logistic Regression.
2. Check if there is class imbalance problem in Dataset 2. Compare the results of decision tree classifier and AdaBoost classifier on Dataset 2 and write your analysis.



### **Department of Computer Science and Engineering (Data Science)**

3. Implement AdaBoost with base learner as decision tree on dataset 2 using K fold cross validation. Perform Hyperparameter tuning using (a) different depth, (b) different learning rate and (c) grid search CV. Show your results using Boxplot.

COLAB LINKS:

<https://colab.research.google.com/drive/1Fx6KO6jRE3rFbiXlBDcmRQE8nQClABi6?usp=sharing>

<https://colab.research.google.com/drive/1aview3517xWjxRTetkBQ6etsCpwl1lB?usp=sharing>