

Statistics Foundation

Machine Learning – Classification

Classification Problem

- **What is a Classification Problem ?**

- We identify problem as classification problem when independent variables are continuous in nature and dependent variable is in categorical form i.e. in classes like positive class and negative class

- The real life example of classification example would be

- to categorize the mail as spam or not spam
- to categorize the tumor as malignant or benign
- to categorize the transaction as fraudulent or genuine

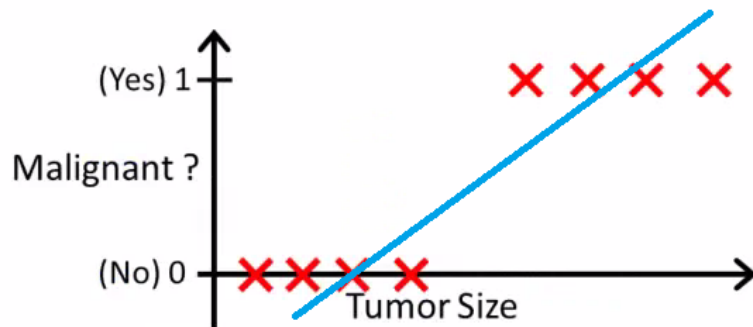
- All these problem's answers are in categorical form i.e. Yes or No and that is why they are two class classification problems

- Although, sometime we come across more than 2 classes and still it is a classification problem. These types of problems are known as multi class classification problems.

Two Class Classification		
$y \in \{0, 1\}$	1 or Positive Class	0 or Negative Class
Email	Spam	Not Spam
Tumor	Malignant	Benign
Transaction	Fraudulent	Not Fraudulent

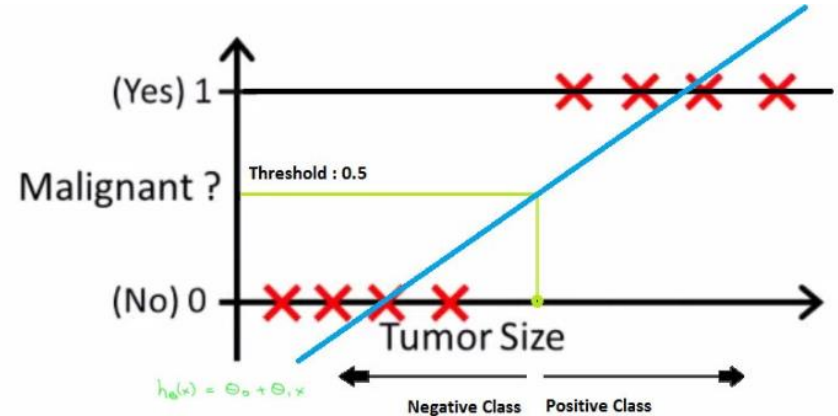
Classification Problem

- **Why not use Linear Regression ?**
 - Suppose we have a data of tumour size vs its malignancy.
 - As it is a classification problem, if we plot, we can see, all the values will lie on 0 and 1.
 - And if we fit best found regression line, by assuming the threshold at 0.5, we can do line pretty reasonable job.



Classification Problem

- **Why not use Linear Regression ?**
 - We can decide the point on the x axis from where:
 - All the values lie to its left side are considered as negative class
 - All the values lie to its right side are positive class



Classification Problem

- **Why not use Linear Regression ?**

- What if there is an outlier in the data
- Things would get pretty messy
- For example, for 0.5 threshold

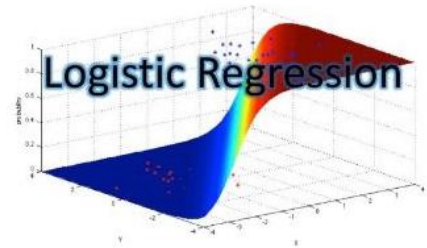


- If we fit best found regression line, it still won't be enough to decide any point by which we can differentiate classes
- It will put some positive class examples into negative class
- The green dotted line (Decision Boundary) is dividing malignant tumours from benign tumours but the line should have been at a yellow line which is clearly dividing the positive and negative examples
- So just a single outlier is disturbing the whole linear regression predictions
- And that is where logistic regression comes into a picture

Logistic Regression

- **What is Logistic Regression**

- Logistic Regression is a Machine Learning classification algorithm that is used to predict the probability of a categorical dependent variable
- In logistic regression, the dependent variable is a binary variable that contains data coded as 1 (yes, success, etc.) or 0 (no, failure, etc.)
- In other words, the logistic regression model predicts $P(Y=1)$ as a function of X
- It is named as 'Logistic Regression', because it's underlying technique is quite the same as Linear Regression. The term "Logistic" is taken from the Logit function that is used in this method of classification



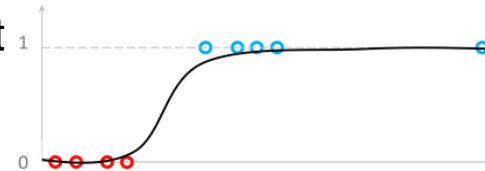
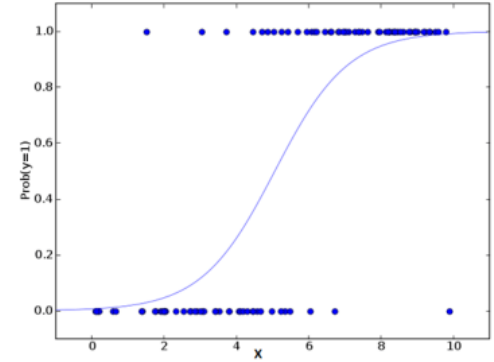
Logistic Regression

- **Logistic Regression**

- Logistic Regression uses Sigmoid function
- The logistic function is a Sigmoid function, which takes any real value between 0 and 1

$$\sigma(t) = \frac{e^t}{e^t + 1} = \frac{1}{1 + e^{-t}}$$

- Let's consider t as linear function in a univariate regression model $t = \beta_0 + \beta_1 x$
- So the Logistic Equation will become $p(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$
- Now, when logistic regression model come across an outlier, it will take care of it



Logistic Regression

- **Logistic Regression Equation**
 - The underlying algorithm of Maximum Likelihood Estimation (MLE) determines the regression coefficient for the model that accurately predicts the probability of the binary dependent variable
 - The algorithm stops when the convergence criterion is met or maximum number of iterations are reached
 - Since the probability of any event lies between 0 and 1 (or 0% to 100%), when we plot the probability of dependent variable by independent factors, it will demonstrate an 'S' shape curve.

$\text{Logit} = \text{Log} (p/1-p) = \log (\text{probability of event happening} / \text{probability of event not happening}) = \log (\text{Odds})$

Logistic Regression

- **Logistic Regression Example**

- We are provided a sample of 1000 customers. We need to predict the probability whether a customer will buy (y) a particular magazine or not. As we've a categorical outcome variable, we'll use logistic regression

- To start with logistic regression, first write the simple linear regression equation with dependent variable enclosed in a link function:

$$g(y) = \beta_0 + \beta(\text{Age})$$

- For understanding, consider 'Age' as independent variable

- In logistic regression, we are only concerned about the probability of outcome dependent variable (success or failure).

- $g()$ is the link function. This function is established using two things:

- Probability of Success(p) and Probability of Failure($1-p$)

- p should meet following criteria:

- It must always be positive (since $p \geq 0$)

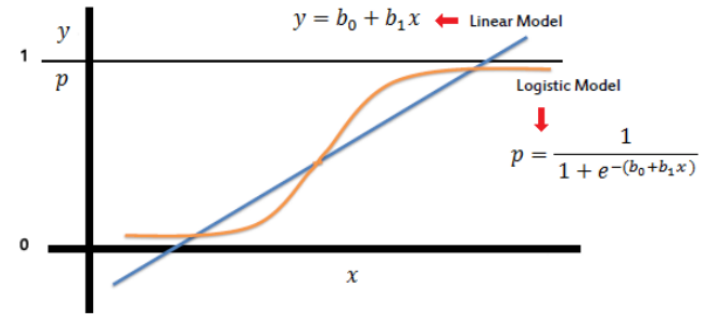
- It must always be less than equals to 1 (since $p \leq 1$)

- Final equation for Logistic Regression:

$$\log(p/(1-p)) = \beta_0 + \beta(\text{Age})$$

Logistic Regression

- **Logistic Regression key facts**
 - $(p/1-p)$ is the odd ratio
 - Whenever the log of odd ratio is found to be positive, the probability of success is always more than 50%
 - A typical logistic model plot is shown below
 - It shows probability never goes below 0 and above 1
 - Logistic regression is similar to a linear regression, but the curve is constructed using the natural logarithm of the “odds” of the target variable, rather than the probability
 - Moreover, the predictors do not have to be normally distributed or have equal variance in each group
 - Logistic regression can handle any number of numerical and/or categorical variables



$$p = \frac{1}{1 + e^{-(b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p)}}$$

Logistic Regression

Logistic Regression Evaluation

- **AIC** (Akaike Information Criteria) – The analogous metric of adjusted R^2 in logistic regression is AIC. AIC is the measure of fit which penalizes model for the number of model coefficients. Therefore, we always prefer model with minimum AIC value

- **Null Deviance and Residual Deviance** – Null Deviance indicates the response predicted by a model with nothing but an intercept. Lower the value, better the model. Residual deviance indicates the response predicted by a model on adding independent variables. Lower the value, better the model

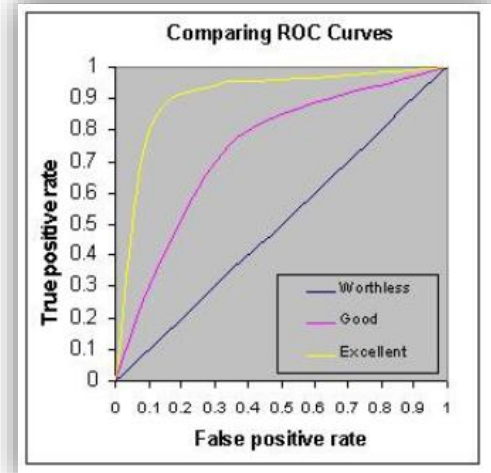
- **Confusion Matrix:** This is a tabular representation of Actual vs Predicted values. This helps us to find the accuracy of the model and avoid over-fitting

		Predicted Class		
		No	Yes	
Observed Class	No	TN	FP	Model Performance
	Yes	FN	TP	
TN	True Negative			Accuracy = $(TN+TP)/(TN+FP+FN+TP)$
FP	False Positive			Precision = $TP/(FP+TP)$
FN	False Negative			Sensitivity = $TP/(TP+FN)$
TP	True Positive			Specificity = $TN/(TN+FP)$

Confusion Matrix

Logistic Regression

- **Logistic Regression Evaluation**
 - **ROC Curve:** Receiver Operating Characteristic (ROC) summarizes the model's performance by evaluating the trade-offs between true positive rate (sensitivity) and false positive rate (1- specificity)
 - For plotting ROC, it is advisable to assume $p > 0.5$ since we are more concerned about success rate
 - ROC summarizes the predictive power for all possible values of $p > 0.5$
 - The area under curve (AUC), referred to as index of accuracy (A) or concordance index, is a perfect performance metric for ROC curve
 - Higher the area under curve, better the prediction power of the model
 - The ROC of a perfect predictive model has TP equals 1 and FP equals 0. This curve will touch the top left corner of the graph



Logistic Regression

- **Probability, Odds and Log of Odds**
 - Let's say that the probability of success is $p=0.8$, then the probability of failure is $1-p=0.2$
 - The odds of success is $p/(1-p)=0.8/0.2=4$, i.e. the odds of success is 4 to 1 and the odds of failure is 0.25 to 1
 - Note that:
 - Probability ranges from 0 to 1
 - Odds range from 0 to ∞
 - Log Odds range from $-\infty$ to ∞

Logistic Regression

- **Interpretation of Logistic Output**
 - Consider sample dataset 'honordata'

##	female	read	write	math	hon	femalexmth
## 1	0	57	52	41	0	0
## 2	1	68	59	53	0	53
## 3	0	44	33	54	0	0
## 4	0	63	44	47	0	0
## 5	0	47	52	57	0	0
## 6	0	44	52	51	0	0

Logistic Regression

- **Interpretation of Logistic Output**

Step1:

- Logistic Regression with No Predictor Variables
- Target variable is 'hon' (as shown sample data)

$$\text{logit}(p) = \beta_0$$

```
##          Estimate Std. Error  z value    Pr(>|z|)
## (Intercept) -1.12546    0.16441 -6.845446 7.623779e-12
```

##	female	read	write	math	hon	female	math
## 1	0	57	52	41	0		0
## 2	1	68	59	53	0		53
## 3	0	44	33	54	0		0
## 4	0	63	44	47	0		0
## 5	0	47	52	57	0		0
## 6	0	44	52	51	0		0

The intercept= -1.12546 which corresponds to the log odds of the probability of being in an honor class p

We can go from the log odds to the odds by exponentiating the coefficient which gives us the odds $O=0.3245$

We can go backwards to the probability by calculating $p=O/(1+O) = 0.245$

Logistic Regression

- **Interpretation of Logistic Output**

Step2:

- Logistic Regression with a Single Dichotomous Predictor Variable
- Target variable is 'hon' (as shown sample data)
- Here we will use a binary predictor variable female in our model:

$$\text{logit}(p) = \beta_0 + \beta_1 * \text{female}$$

```
##           Estimate Std. Error  z value    Pr(>|z|)
## (Intercept) -1.4708517  0.2689554 -5.468756 4.532047e-08
## female      0.5927822  0.3414293  1.736178 8.253231e-02
```

##	female	read	write	math	hon	femalexmath
## 1	0	57	52	41	0	0
## 2	1	68	59	53	0	53
## 3	0	44	33	54	0	0
## 4	0	63	44	47	0	0
## 5	0	47	52	57	0	0
## 6	0	44	52	51	0	0

The intercept= -1.47085 which corresponds to the log odds for males being in an honor class (since male is the reference group, female=0)

The coefficient for female= 0.59278 which corresponds to the log of odds ratio between the female group and male group
The odds ratio equals 1.81 which means the odds for females are about 81% higher than the odds for males

Logistic Regression

- **Interpretation of Logistic Output**

Step3:

- Logistic Regression with a Single Continuous Predictor Variable
- Target variable is 'hon' (as shown sample data)
- Here we will use a single continuous predictor variable math in our model:

$$\text{logit}(p) = \beta_0 + \beta_1 * \text{math}$$

```
##           Estimate Std. Error  z value    Pr(>|z|)
## (Intercept) -9.7939421  1.48174484 -6.609736 3.850061e-11
## math        0.1563404  0.02560948  6.104784 1.029399e-09
```

##	female	read	write	math	hon	female	math
## 1	0	57	52	41	0		0
## 2	1	68	59	53	0		53
## 3	0	44	33	54	0		0
## 4	0	63	44	47	0		0
## 5	0	47	52	57	0		0
## 6	0	44	52	51	0		0

The intercept= -9.79394 which is interpreted as the log odds of a student with a math score of zero being in an honors class

The coefficient for math= 0.15634 which is interpreted as the expected change in log odds for a one-unit increase in math score
The odds ratio can be calculated by exponentiating this value to get 1.16922 which means we expect to see about 17% increase in the odds of being in an honors class, for a one-unit increase in math score

$$\text{logit}(p) = \frac{p}{1-p} = -9.79394 + 0.15634 * \text{math}$$

Logistic Regression

- **Python Implementation**

```
import pandas as pd
import numpy as np
Diabetes=pd.read_csv('diabetes.csv')
table1=np.mean(Diabetes,axis=0)
table2=np.std(Diabetes,axis=0)
```

```
#####The data are unbalanced with 35% of observations having diabetes.
#####The standard deviation of the different variables is also very
different, to compare the coefficient of the different variables the
coefficient will need to be standardized
```

```
inputData=Diabetes.iloc[:,8]
outputData=Diabetes.iloc[:,8]
```

Logistic Regression

- **Python Implementation**

```
from sklearn.linear_model import LogisticRegression  
logit1=LogisticRegression()  
logit1.fit(inputData,outputData)
```

```
logit1.score(inputData,outputData)
```

Even if the logistic regression is a simple model around 78% of the observation are correctly classified!

Due to class imbalance, we need to check the model performance on each class. Not being able to classify people with diabetes would be a major problem since this is the goal of the model.

Logistic Regression

- **Python Implementation**

```
##True positive
trueInput=Diabetes.ix[Diabetes['Outcome']==1].iloc[:,8]
trueOutput=Diabetes.ix[Diabetes['Outcome']==1].iloc[:,8]
##True positive rate
np.mean(logit1.predict(trueInput)==trueOutput)
##Return around 55%

##True negative
falseInput=Diabetes.ix[Diabetes['Outcome']==0].iloc[:,8]
falseOutput=Diabetes.ix[Diabetes['Outcome']==0].iloc[:,8]
##True negative rate
np.mean(logit1.predict(falseInput)==falseOutput)
##Return around 90%
```

Logistic Regression

- **Python Implementation**

```
###Confusion matrix with sklearn  
from sklearn.metrics import confusion_matrix, roc_curve, roc_auc_score  
confusion_matrix(logit1.predict(inputData),outputData)
```

Logistic Regression

- **Python Implementation**

##Computing false and true positive rates

```
fpr, tpr, _ = roc_curve(logit1.predict(inputData), outputData, drop_intermediate=False)
```

```
import matplotlib.pyplot as plt
```

```
plt.figure()
```

```
##Adding the ROC
```

```
plt.plot(fpr, tpr, color='red', lw=2, label='ROC curve')
```

```
##Random FPR and TPR
```

```
plt.plot([0, 1], [0, 1], color='blue', lw=2, linestyle='--')
```

```
##Title and label
```

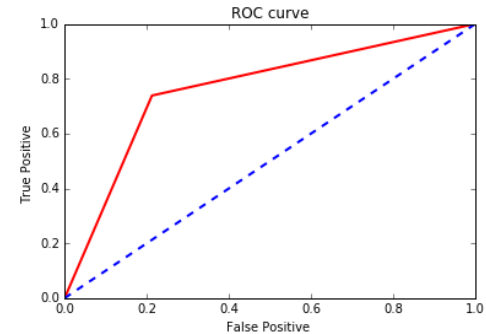
```
plt.xlabel('FPR')
```

```
plt.ylabel('TPR')
```

```
plt.title('ROC curve')
```

```
plt.show()
```

```
roc_auc_score(logit1.predict(inputData), outputData)
```



The ROC curve of the model

Decision Tree

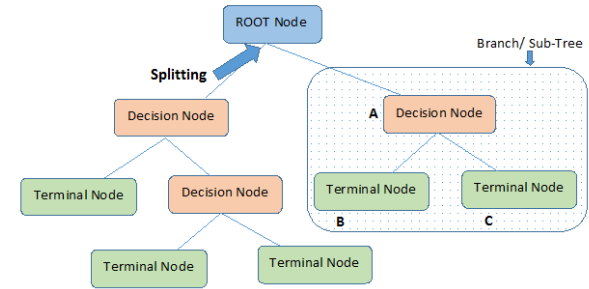
- **Why Decision Tree**
 - They are supervised learning algorithm which has a pre-defined target variable
 - They are mostly used in non-linear decision making with simple linear decision surface
 - They are adaptable for solving both classification or regression kind of problems

Decision Tree

- **Why Decision Tree**
 - They empower predictive modeling with higher accuracy, better stability and provide ease of interpretation.
 - Unlike linear modeling techniques, they map non-linear relationships quite well.

Decision Tree

- **How Decision Tree works**
 - It breaks a dataset into smaller subsets, and at the same time, an associated decision tree is incrementally developed.
 - As it happens in real life, we consider the most important factor, and divide possibilities according to it.
 - Similarly, tree building starts by finding the variable/feature for best split.

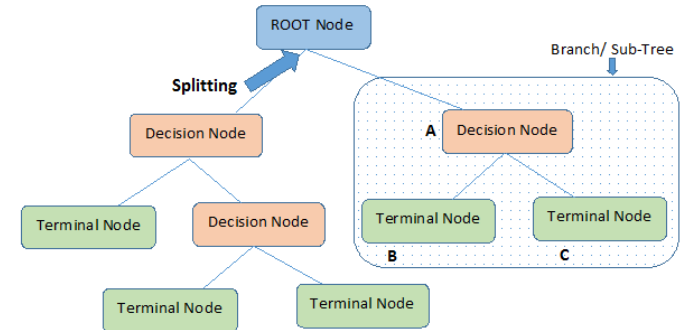


Note:- A is parent node of B and C.

Terminology explained.

Decision Tree

- **Decision Tree Terminologies**
 - Root Node: Entire population or sample, further gets divided into two or more homogeneous sets.
 - Parent and Child Node: Node which is divided into sub-nodes is called parent node, whereas sub-nodes are the child of parent node.
 - Splitting: Process of dividing a node into two or more sub-nodes.
 - Decision Node: A sub-node that splits into further sub-nodes.
 - Leaf/ Terminal Node: Nodes that do not split.
 - Pruning: When we remove sub-nodes of a decision node, this process is called pruning. (Opposite of Splitting)
 - Branch/Sub-Tree: Sub-section of entire tree.



Note:- A is parent node of B and C.

Terminology explained.

Decision Tree

- **What is a Decision Tree**
 - A decision tree consists of the root /Internal node which further splits into decision nodes/branches, depending on the outcome of the branches the next branch or the terminal /leaf nodes are formed

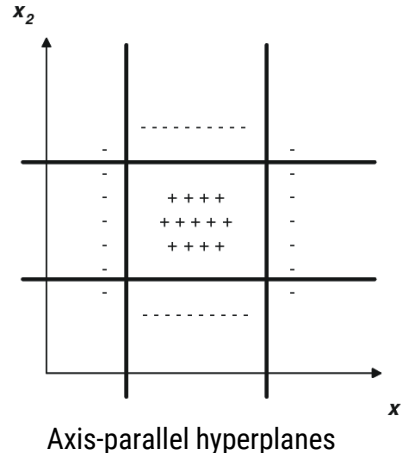
Alternatingly think that decision trees are a group of nested IF-ELSE conditions which can be modeled as a tree wherein the decision are made in the internal node and output is obtained in the leaf node.

```
if (Height > 5.6):  
    output = TOUCHDOWN  
elif (speed < 13.2 m.p.h):  
    output = INTERCEPTION  
elif (weight > 200 lbs):  
    output = TOUCHDOWN  
else:  
    output = INTERCEPTION
```

Nested if-else

Decision Tree

- **Geometrical analogy:**
 - Geometrically we can think decision trees as a set of a number of axis parallel hyperplanes which divides the space into the number of hypercuboids during the inference. we classify the point based on which hypercuboid it falls into.



Decision Tree

- **What is Entropy method ?**
 - Definition: Entropy is the measures of impurity, disorder or uncertainty in a bunch of examples.
 - Entropy controls how a Decision Tree decides to split the data. It actually effects how a Decision Tree draws its boundaries.

The Equation of Entropy:

$$\text{Entropy} = - \sum p(X) \log p(X)$$



here $p(x)$ is a fraction of
examples in a given class

Equation of Entropy

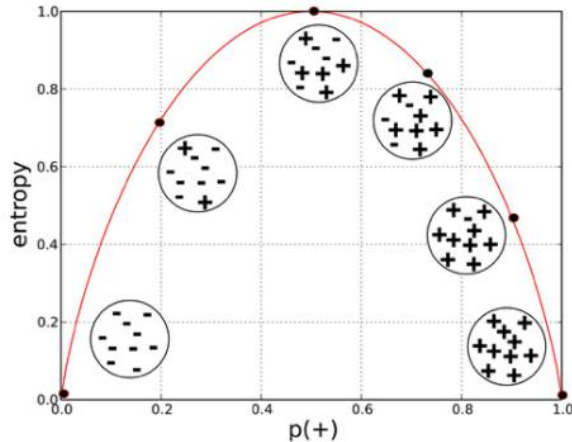
Decision Tree

- **What is Entropy method ?**
 - Let's say we only have two classes , a positive class and a negative class. Therefore 'i' here could be either + or (-).
 - So if we had a total of 100 data points in our dataset with 30 belonging to the positive class and 70 belonging to the negative class
 - Then 'P+' would be 3/10 and 'P-' would be 7/10

$$-\frac{3}{10} \times \log_2\left(\frac{3}{10}\right) - \frac{7}{10} \times \log_2\left(\frac{7}{10}\right) \approx 0.88$$

Entropy in this example

Decision Tree



Entropy variation for Pure/Impure Nodes

What is Entropy method ?

- The x-axis measures the proportion of data points belonging to the positive class in each bubble and the y-axis axis measures their respective entropies.
- We can see the inverted 'U' shape of the graph. Entropy is lowest at the extremes, when the bubble either contains no positive instances or only positive instances.
- That is, when the bubble is pure the disorder is 0. Entropy is highest in the middle when the bubble is evenly split between positive and negative instances. Extreme disorder , because there is no majority.

Decision Tree

$$IG(Y, X) = E(Y) - E(Y|X)$$

Information Gain from X on Y

- **What is Information Gain ?**
 - Now we know how to measure disorder (using Entropy).
 - Next we need a metric to measure the reduction of this disorder in our target variable/class given additional information(features/independent variables) about it. This is where Information Gain comes in

Decision Tree

$$IG(Y, X) = E(Y) - E(Y|X)$$

Information Gain from X on Y

- **What is Information Gain ?**
 - We simply subtract the entropy of Y given X from the entropy of just Y to calculate the reduction of uncertainty about Y given an additional piece of information X about Y
 - This is called Information Gain
 - The greater the reduction in this uncertainty, the more information is gained about Y from X.

Decision Tree

- **Example: Decision Tree**

Consider an example where we are building a decision tree to predict whether a loan given to a person would result in a write-off or not :

- Our entire population consists of 30 instances
- 16 belong to the write-off class
- Other 14 belong to the non-write-off class
- We have two features, namely “Balance” that can take on two values -> “< 50K” or “>50K”
- “Residence” that can take on three values -> “OWN”, “RENT” or “OTHER”

Let us see how a decision tree algorithm would decide what attribute to split on first and what feature provides more information, or reduces more uncertainty about our target variable out of the two using the concepts of Entropy and Information Gain.

Decision Tree

- **Example: Decision Tree**

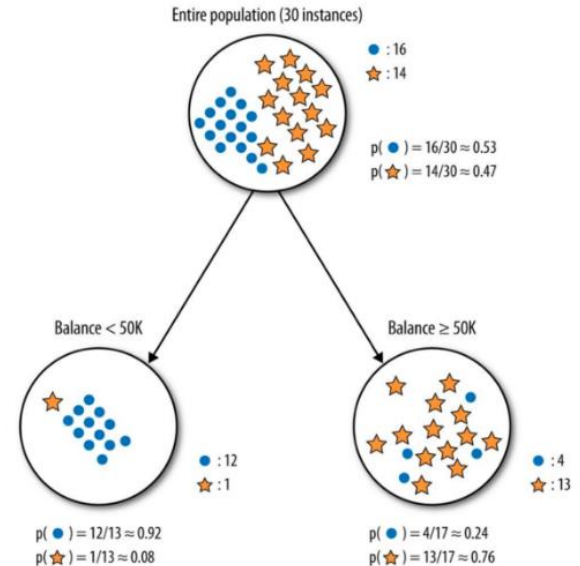
Split on Feature 1: Balance

The dots are the data points with class right-off and the stars are the non-write-offs

Splitting the parent node on attribute balance gives us 2 child nodes

The left node gets 13 of the total observations with 12/13 (0.92 probability) observations from the write-off class and only 1/13(0.08 probability) observations from the non-write of class

The right node gets 17 of the total observation with 13/17(0.76 probability) observations from the non-write-off class and 4/17 (0.24 probability) from the write-off class



Decision Tree

Example: Decision Tree

Split on Feature 1: Balance

Let's calculate the entropy for the parent node and see how much uncertainty the tree can reduce by splitting on Balance

$$E(\text{Parent}) = -\frac{16}{30}\log_2\left(\frac{16}{30}\right) - \frac{14}{30}\log_2\left(\frac{14}{30}\right) \approx 0.99$$

$$E(\text{Balance} < 50K) = -\frac{12}{13}\log_2\left(\frac{12}{13}\right) - \frac{1}{13}\log_2\left(\frac{1}{13}\right) \approx 0.39$$

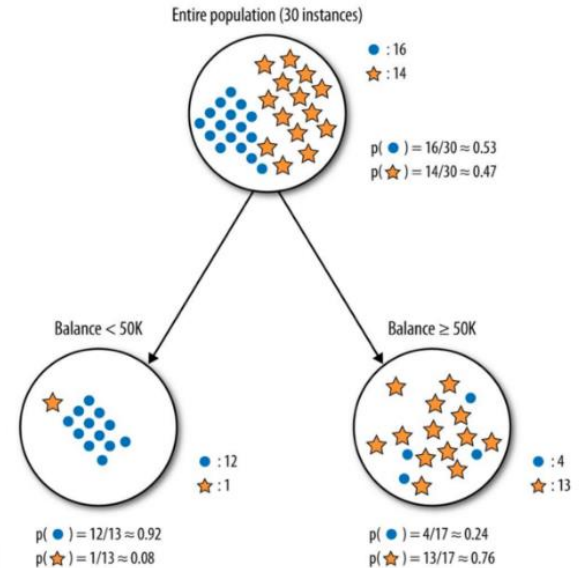
$$E(\text{Balance} > 50K) = -\frac{4}{17}\log_2\left(\frac{4}{17}\right) - \frac{13}{17}\log_2\left(\frac{13}{17}\right) \approx 0.79$$

Weighted Average of entropy for each node:

$$\begin{aligned} E(\text{Balance}) &= \frac{13}{30} \times 0.39 + \frac{17}{30} \times 0.79 \\ &= 0.62 \end{aligned}$$

Information Gain:

$$\begin{aligned} IG(\text{Parent}, \text{Balance}) &= E(\text{Parent}) - E(\text{Balance}) \\ &= 0.99 - 0.62 \\ &= 0.37 \end{aligned}$$



Decision Tree

Example: Decision Tree

Split on Feature 2: Residence

Let's calculate the entropy for the parent node and see how much uncertainty the tree can reduce by splitting on Residence

$$E(\text{Residence} = \text{OWN}) = -\frac{7}{8}\log_2\left(\frac{7}{8}\right) - \frac{1}{8}\log_2\left(\frac{1}{8}\right) \approx 0.54$$

$$E(\text{Residence} = \text{RENT}) = -\frac{4}{10}\log_2\left(\frac{4}{10}\right) - \frac{6}{10}\log_2\left(\frac{6}{10}\right) \approx 0.97$$

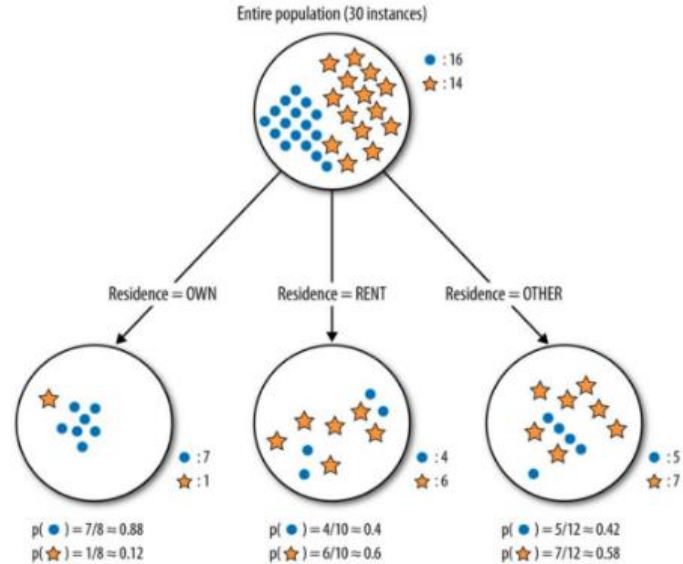
$$E(\text{Residence} = \text{OTHER}) = -\frac{5}{12}\log_2\left(\frac{5}{12}\right) - \frac{7}{12}\log_2\left(\frac{7}{12}\right) \approx 0.98$$

Weighted Average of entropies for each node:

$$E(\text{Residence}) = \frac{8}{30} \times 0.54 + \frac{10}{30} \times 0.97 + \frac{12}{30} \times 0.98 = 0.86$$

Information Gain:

$$\begin{aligned} IG(\text{Parent}, \text{Residence}) &= E(\text{Parent}) - E(\text{Residence}) \\ &= 0.99 - 0.86 \\ &= 0.13 \end{aligned}$$



Decision Tree

- By itself the feature, Balance provides more information about our target variable than Residence
- It reduces more disorder in our target variable.
- A decision tree algorithm would use this result to make the first split on our data using Balance
- From here on, the decision tree algorithm would use this process at every split to decide what feature it is going to split on next
- In a real world scenario , with more than two features the first split is made on the most informative feature and then at every split the information gain for each additional feature needs to be done
- A decision tree would repeat this process as it grows deeper and deeper till either it reaches a pre-defined depth or no additional split can result in a higher information gain beyond a certain threshold which can also usually be specified as a hyper-parameter

Decision Tree

- **What is Gini Impurity**

Gini impurity can be considered as an alternative for the entropy method. Gini impurity is a measure of how often a randomly chosen element from the set would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset.

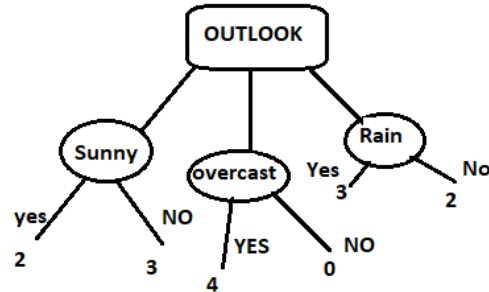
$$G.I(Y) = 1 - \sum_{i=1}^k [p(Y)]^2$$

Calculation of Entropy in Gini method behaves in the same way but entropy involves a log computation and Gini impurity involved a square computation. Since computing square is cheaper than logarithmic function we prefer Gini impurity over entropy.

Decision Tree

- **What is Pure Node**

Pure node is a node wherein all the data points belong to the same class and thus it is very easy to make the prediction at such node.



Here overcast is a pure node

Decision Tree

- **Overfitting & Underfitting in CART**
- If we continue to grow the tree fully until each leaf node corresponds to the lowest impurity, then the data have typically been overfitted
- If splitting is stopped too early, error on training data is not sufficiently high and performance will suffer due to bias
- Thus, preventing overfitting & underfitting are pivotal while modeling a decision tree and it can be done in 2 ways:
 - **Setting constraints on tree size**
 - **Tree pruning**

Decision Tree

- **Overfitting & Underfitting in CART**
 - **Setting constraints on tree size**
 - Providing minimum number of samples for a node split.
 - Deploying the minimum number of samples for a terminal node (leaf).
 - Allowing maximum depth of tree (vertical depth).
 - Maximum number of terminal nodes.
 - Maximum features to consider for the split.

Decision Tree

- **Overfitting & Underfitting in CART**
 - **Tree Pruning**
 - Pruning is a technique in machine learning that reduces the size of decision trees by removing sections of the tree
 - It also reduces the complexity of the final classifier, and hence improves predictive accuracy by the reduction of overfitting
 - Tree pruning can be done in two ways by pre-pruning or by post-pruning

Decision Tree

- **Overfitting & Under-fitting in CART**
 - **Pre-pruning:**
 - Stop splitting the current node if it does not improve the entropy by at least some pre-set(threshold) value
 - Stop partitioning if the number of data points are less then some pre-set (Threshold) values
 - Restricting the depth of the tree to some pre-set(Threshold) value

- **Post-pruning:**

It can be done by first allowing the tree to grow to its full potential and then pruning the tree at each level after calculating the cross-validation accuracy at each level.

Decision Tree

- **Advantages of CART:**
 - Decision trees can inherently perform multiclass classification
 - They provide most model interpretability because they are simply series of if-else conditions
 - They can handle both numerical and categorical data
 - Nonlinear relationships among features do not affect the performance of the decision trees

Decision Tree

- **Disadvantages of CART:**
 - A small change in the dataset can make the tree structure unstable which can cause high variance
 - Decision tree learners create underfit trees if some classes are imbalanced. It is therefore recommended to balance the data set prior to fitting with the decision tree

Decision Tree

- **Preparing data for CART:**
 - The splitting of numerical features can be performed by sorting the features in the ascending order and trying each value as the threshold point and calculating the information gain for each value as the threshold. Finally, if that value obtained is equal to the threshold which gives the maximum I.G value then that point is chosen as the split point
 - Feature scaling(column standardization) not necessary to perform in decision trees. However, it helps with data visualization/manipulation and might be useful if you intend to compare performance with other data or other methods like SVM.
 - In order to handle categorical features in Decision trees, we must never perform one hot encoding on a categorical variable even if the categorical variables are nominal since most of the libraries can handle categorical variables automatically. we can still assign a number for each variable if desired.
 - If height or depth of the tree is exactly one then such a tree is called as a decision stump.

Decision Tree

- **Preparing data for CART:**
 - Imbalanced class does have a detrimental impact on the tree's structure so it can be avoided by either using up-sampling or by using down-sampling depending upon the dataset
 - Apart from skewed classes, high dimensionality can also have an adverse effect on the structure of the tree if dimensionality is very high that means we have a lot of features which means that to find the splitting criterion on each node it will consume a lot of time
 - Outliers also impact the tree's structure as the depth increases the chance of outliers in the tree increases

Decision Tree

Python Code

```
#Import Library
#Import other necessary libraries like pandas, numpy...
from sklearn import tree
#Assumed you have, X (predictor) and Y (target) for training data set and x_test(predictor)
of test_dataset
# Create tree object
model = tree.DecisionTreeClassifier(criterion='gini') # for classification, here you can
change the algorithm as gini or entropy (information gain) by default it is gini
# model = tree.DecisionTreeRegressor() for regression
# Train the model using the training sets and check score
model.fit(X, y)
model.score(X, y)
#Predict Output
predicted= model.predict(x_test)
```

Ensemble Methods

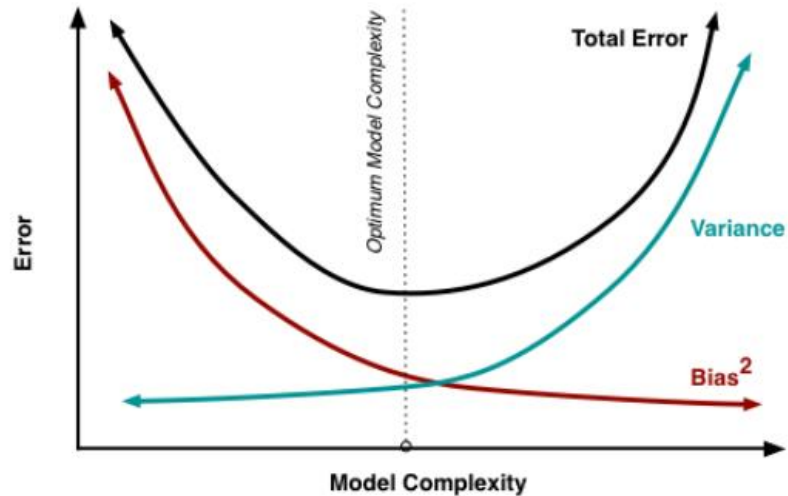
Ensemble Methods

- The literary meaning of word 'ensemble' is group
- Ensemble methods involve group of predictive models to achieve a better accuracy and model stability
- Ensemble methods are known to impart supreme boost to tree based models
- Like every other model, a tree based model also suffers from the plague of bias and variance.
- Bias means, 'how much on an average are the predicted values different from the actual value.'
- Variance means, 'how different will the predictions of the model be at the same point if different samples are taken from the same population'.
- With a small tree we will get a model with low variance and high bias
- How do we balance the trade off between bias and variance ?
- Normally, as the complexity of model increases, there is a reduction in prediction error due to lower bias in the model
- As the model grows more complex, we end up over-fitting model and the model will start suffering from high variance

Ensemble Methods

Ensemble Methods

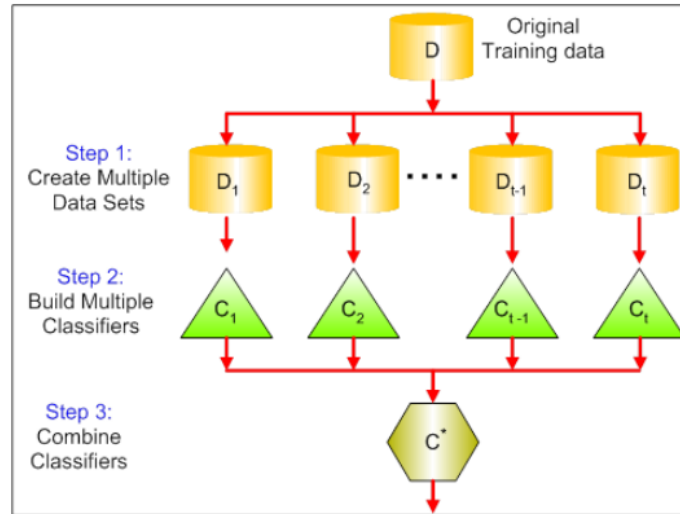
A champion model should maintain a balance between these two types of errors. This is known as the trade-off management of bias-variance errors. Ensemble learning is one way to execute this trade off analysis.



Ensemble Methods

Bagging

Bagging is a technique used to reduce the variance of our predictions by combining the result of multiple classifiers modelled on different sub-samples of the same data set.



Ensemble Methods

Steps to perform Bagging

1. Create Multiple DataSets:

- Sampling is done with replacement on the original data and new datasets are formed.
- The new data sets can have a fraction of the columns as well as rows, which are generally hyper-parameters in a bagging model
- Taking row and column fractions less than 1 helps in making robust models, less prone to overfitting

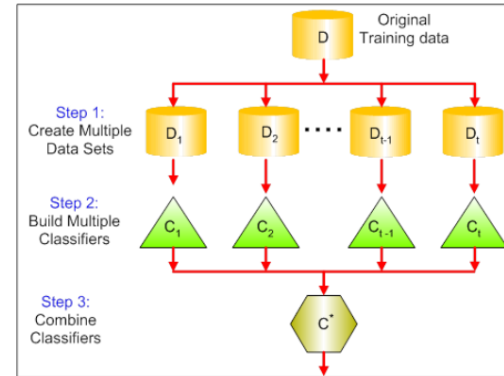
2. Build Multiple Classifiers:

- Classifiers are built on each data set.
- Generally the same classifier is modelled on each data set and predictions are made.

3. Combine Classifiers:

- The predictions of all the classifiers are combined using a mean, median or mode value depending on the problem at hand.
- The combined values are generally more robust than a single model

There are various implementations of bagging models. Random forest is one of them



Random Forests

- **Random Forests**

Random Forest is a versatile machine learning method capable of performing both regression and classification tasks. It also undertakes dimensional reduction methods, treats missing values, outlier values and other essential steps of data exploration, and does a fairly good job. It is a type of ensemble learning method, where a group of weak models combine to form a powerful model.

- In Random Forest, we grow multiple trees as opposed to a single tree in CART model
- To classify a new object based on attributes, each tree gives a classification and we say the tree “votes” for that class
- The forest chooses the classification having the most votes (over all the trees in the forest) and in case of regression, it takes the average of outputs by different trees.

Random Forests

Steps to build Random Forest

- Assume number of cases in the training set is N . Then, sample of these N cases is taken at random but with replacement. This sample will be the training set for growing the tree
- If there are M input variables, a number $m < M$ is specified such that at each node, m variables are selected at random out of the M .
- The best split on these m is used to split the node. The value of m is held constant while we grow the forest
- Each tree is grown to the largest extent possible and there is no pruning
- Predict new data by aggregating the predictions of the n trees (i.e., majority votes for classification, average for regression)



Random Forests

- **Advantages**
 - This algorithm can solve both type of problems i.e. classification and regression and does a decent estimation at both fronts
 - One of benefits of Random forest is, the power of handle large data set with higher dimensionality
 - It can handle thousands of input variables and identify most significant variables so it is considered as one of the dimensionality reduction methods
 - Further, the model outputs Importance of variable, which can be a very handy feature (on some random data set).
 - It has an effective method for estimating missing data and maintains accuracy when a large proportion of the data are missing.

Random Forests

- **Advantages**
 - It has methods for balancing errors in data sets where classes are imbalanced.
 - The capabilities of the above can be extended to unlabelled data, leading to unsupervised clustering, data views and outlier detection
 - Random Forest involves sampling of the input data with replacement called as bootstrap sampling. Here one third of the data is not used for training and can be used to testing. These are called the out of bag samples
 - Error estimated on these out of bag samples is known as out of bag error. Study of error estimates by Out of bag, gives evidence to show that the out-of-bag estimate is as accurate as using a test set of the same size as the training set
 - Therefore, using the out-of-bag error estimate removes the need for a set aside test set

Random Forests

- **Disadvantages**
 - It surely does a good job at classification but not as good as for regression problem as it does not give precise continuous nature predictions.
 - In case of regression, it doesn't predict beyond the range in the training data, and that they may over-fit data sets that are particularly noisy.
 - Random Forest can feel like a black box approach for statistical modellers – we have very little control on what the model does. We can at best – try different parameters and random seeds

Random Forest

Python Code

```
#Import Library
from sklearn.ensemble import RandomForestClassifier #use RandomForestRegressor for
regression problem
#Assumed you have, X (predictor) and Y (target) for training data set and x_test(predictor)
of test_dataset
# Create Random Forest object
model= RandomForestClassifier(n_estimators=1000)
# Train the model using the training sets and check score
model.fit(X, y)
#Predict Output
predicted= model.predict(x_test)
```

Ensemble Method - Boosting

- **Boosting**

Definition: The term 'Boosting' refers to a family of algorithms which converts weak learner to strong learners.

Let's solve a problem of spam email identification:

How would you classify an email as SPAM or not? Like everyone else, our initial approach would be to identify 'spam' and 'not spam' emails using following criteria. If:

- Email has only one image file (promotional image), It's a SPAM
- Email has only link(s), It's a SPAM
- Email body consist of sentence like "You won a prize money of \$ xxxxxx", It's a SPAM
- Email from our official domain "Analyticsvidhya.com", Not a SPAM
- Email from known source, Not a SPAM

Above, we've defined multiple rules to classify an email into 'spam' or 'not spam'. But, do you think these rules individually are strong enough to successfully classify an email? No.

Ensemble Method - Boosting

- **Boosting**

Individually, these rules are not powerful enough to classify an email into 'spam' or 'not spam'.

Therefore, these rules are called as weak learner.

To convert weak learner to strong learner, we'll combine the prediction of each weak learner using methods like:

- Using average/ weighted average
- Considering prediction has higher vote

For example: Above, we have defined 5 weak learners.

Out of these 5, 3 are voted as 'SPAM' and 2 are voted as 'Not a SPAM'.

In this case, by default, we'll consider an email as SPAM because we have higher(3) vote for 'SPAM'.

Ensemble Method - Boosting



Boosting

- An initial model F_0 is defined to predict the target variable y
- This model will be associated with a residual $(y - F_0)$
- A new model h_1 is fit to the residuals from the previous step
- Now, F_0 and h_1 are combined to give F_1 , the boosted version of F_0 . The mean squared error from F_1 will be lower than that from F_0 :

$$F_1(x) \leftarrow F_0(x) + h_1(x)$$

- To improve the performance of F_1 , we could model after the residuals of F_1 and create a new model F_2 :

$$F_2(x) \leftarrow F_1(x) + h_2(x)$$

- This can be done for 'm' iterations, until residuals have been minimized as much as possible:

$$F_m(x) \leftarrow F_{m-1}(x) + h_m(x)$$

Here, the additive learners do not disturb the functions created in the previous steps. Instead, they impart information of their own to bring down the errors.

Gradient Boosting

Python Code

```
#import libraries
from sklearn.ensemble import GradientBoostingClassifier #For Classification
from sklearn.ensemble import GradientBoostingRegressor #For Regression
#use GBM function
clf = GradientBoostingClassifier(n_estimators=100, learning_rate=1.0, max_depth=1)
clf.fit(X_train, y_train)
```

XGBoost

Python Code

```
# First XGBoost model for Pima Indians dataset
from numpy import loadtxt
from xgboost import XGBClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# load data
dataset = loadtxt('pima-indians-diabetes.csv', delimiter=",")

# split data into X and y
X = dataset[:,0:8]
Y = dataset[:,8]

# split data into train and test sets
seed = 7
test_size = 0.33
```

Python Code

```
X_train, X_test, y_train, y_test = train_test_split(X, Y,
test_size=test_size, random_state=seed)

# fit model no training data
model = XGBClassifier()
model.fit(X_train, y_train)

# make predictions for test data
y_pred = model.predict(X_test)
predictions = [round(value) for value in y_pred]

# evaluate predictions
accuracy = accuracy_score(y_test, predictions)
print("Accuracy: %.2f%%" % (accuracy * 100.0))
```


XGBoost Features

- **Unique features of XGBoost**
 - XGBoost is a popular implementation of gradient boosting. Let's discuss some features of XGBoost that make it so interesting
 - Regularization: XGBoost has an option to penalize complex models through both L1 and L2 regularization. Regularization helps in preventing overfitting
 - Handling sparse data: Missing values or data processing steps like one-hot encoding make data sparse. XGBoost incorporates a sparsity-aware split finding algorithm to handle different types of sparsity patterns in the data
 - Weighted quantile sketch: Most existing tree based algorithms can find the split points when the data points are of equal weights (using quantile sketch algorithm). However, they are not equipped to handle weighted data. XGBoost has a distributed weighted quantile sketch algorithm to effectively handle weighted data

XGBoost Features

- **Unique features of XGBoost**
 - Block structure for parallel learning: For faster computing, XGBoost can make use of multiple cores on the CPU. This is possible because of a block structure in its system design. Data is sorted and stored in in-memory units called blocks. Unlike other algorithms, this enables the data layout to be reused by subsequent iterations, instead of computing it again. This feature also serves useful for steps like split finding and column sub-sampling
 - Cache awareness: In XGBoost, non-continuous memory access is required to get the gradient statistics by row index. Hence, XGBoost has been designed to make optimal use of hardware. This is done by allocating internal buffers in each thread, where the gradient statistics can be stored
 - Out-of-core computing: This feature optimizes the available disk space and maximizes its usage when handling huge datasets that do not fit into memory

Thank You