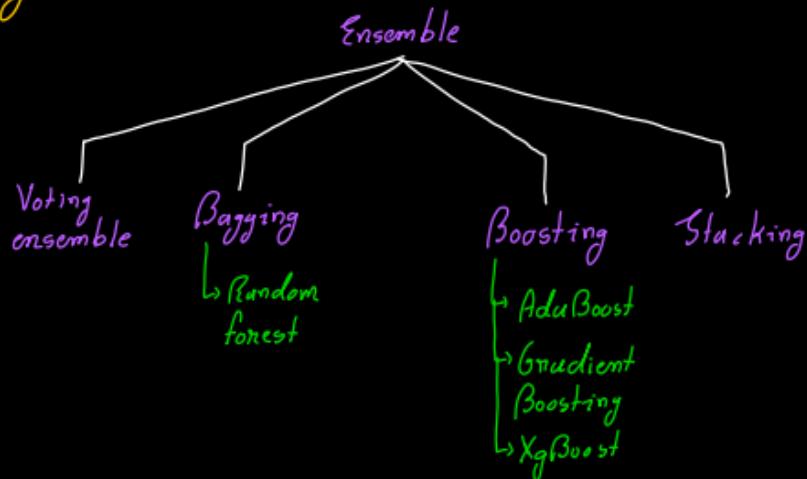


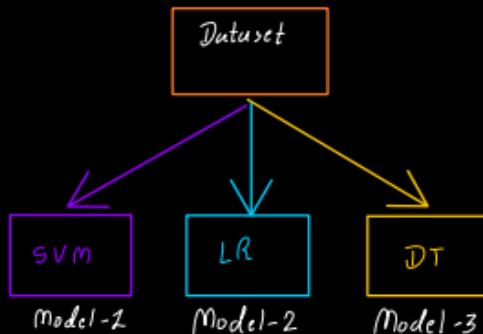
Ensemble learning

→ Many are smarter than few.



1. Voting Ensemble

→ Base models are different algorithms



→ For classification → Majority count

→ For regression → Average



AdaBoost Classification

1. Weak learners

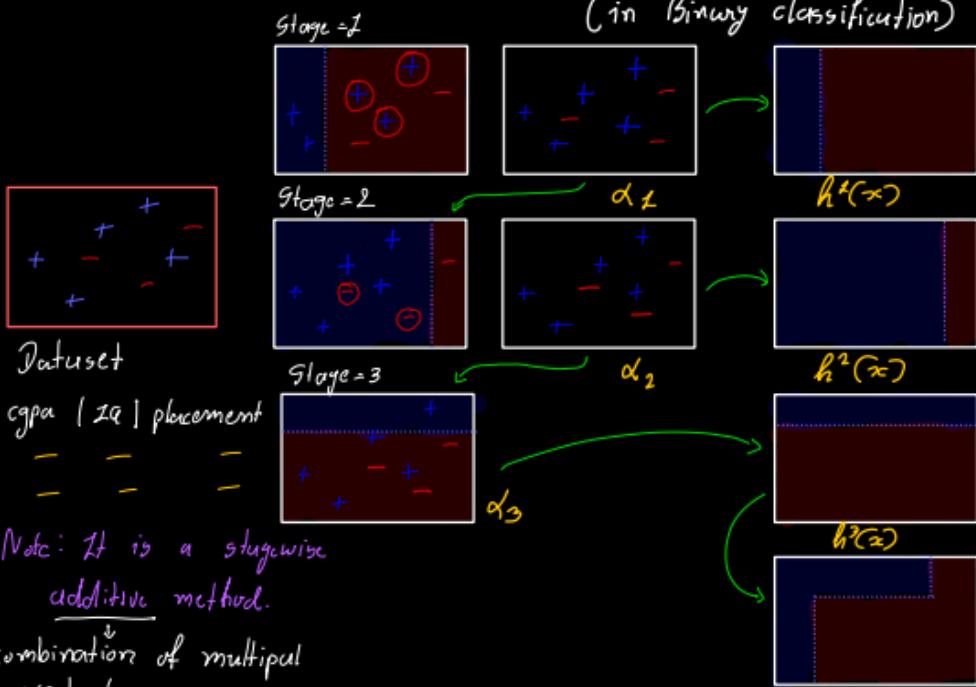
↳ An ML model whose accuracy is very weak (just around 50%).

2. Decision Stumps

↳ A decision tree whose max-depth = 1

3. $+1/-1$

↳ Output will be $+1(+ve)$ or $-1(-ve)$ instead of 0 or 1
(in Binary classification)



combination of multiple weak learners.

→ There can be more number of stages.

→ α_i : weight assigned to each model based on their performance.

→ Now combine all stages

$$H(x) = \text{Sign}(\alpha_1 * h^1(x) + \alpha_2 * h^2(x) + \alpha_3 * h^3(x))$$

→ Calculate sign of entire equation, hence +ve output will be converted to '1' & -ve will be '-1'.
 → Initially 'd' will be assigned as $\frac{1}{n}$ where n = num. of samples.

x_1	x_2	y	\hat{y}	weights
3	7	1	1	0.2
2	9	0	0	0.2
1	4	1	0	0.2
9	8	0	0	0.2
3	7	0	0	0.2

$$n=5, \text{ weights} = \frac{1}{n} = \frac{1}{5} = 0.2$$

Stage = 1

→ Train a Decision tree on training data, whose max-depth = 1.

creating decision stumps like....

$$1. \boxed{x_2 \geq 5} \quad 2. \boxed{x_2 \leq 20}$$

→ Select those stumps, by which entropy decreases maximum or maximum information gain.

suppose we have selected 1st decision stump as M1 which gives

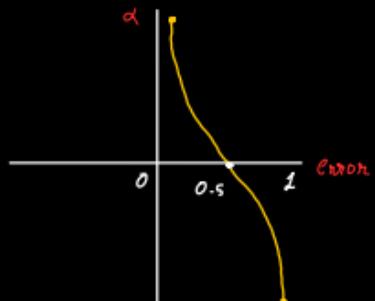
maximum information.

→ Calculate α_1 (for M1) which will depend on your "error rate"

• error rate \uparrow $\alpha_1 \downarrow$ • error rate \downarrow $\alpha_1 \uparrow$

$$\alpha_1 = \frac{1}{2} \ln \left(\frac{1-\text{error}}{\text{error}} \right)$$

Error = Sum of weights of all mis-classified rows.



$$\alpha_1 = \frac{1}{2} \ln \left(\frac{2-\text{error}}{\text{error}} \right)$$

$$\frac{1}{2} \ln \left(\frac{1-0.2}{0.2} \right)$$

$$\frac{1}{2} \ln \left(\frac{0.8}{0.2} \right)$$

$$\alpha_1 = \boxed{0.7}$$

Up-sampling

→ Increasing the weight of wrongly classified rows

→ Decreasing the weight of correctly classified rows

$$\text{new_weight} = \text{current_weight} * e^{\alpha_1}$$

$$\text{new_weight} = \text{current_weight} \boxed{0.7}$$

$$W_{\text{correct}} = 0.2 * e^{-0.7} = \boxed{0.0994}$$

$$W_{\text{wrong}} = 0.2 * e^{0.7} = \boxed{0.4026}$$

Gradient Boosting

- Process of making a powerful ML model by combining multiple weak models.
- Process goes in sequence, where first model will be trained on data.
- The task of next model is to correct the mistakes made by previous model.
- If we are working with regression problem in gradient boosting then the o/p of first model will always be the mean of output column.
- Actually the first node is not a ML model. It is just a mean calculator which always returns mean of output variable. No matter what input is given.

$$\text{Cost} = y - \hat{y} = \text{pseudo residual}$$

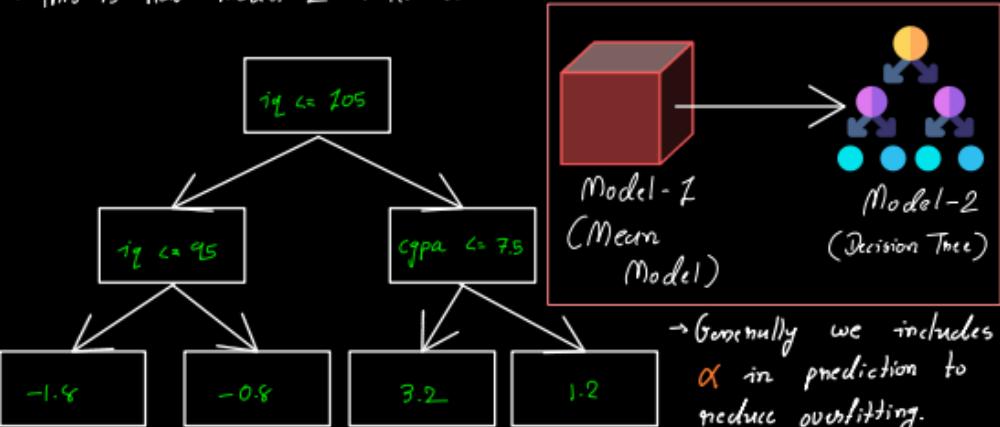
ZQ	CGPA	Salary	Prediction 2(m1)	Pseudo residual	Prediction 2(m2)
90	5	3	4.8	-1.8	-2.8
100	7	4	4.8	-0.8	-0.8
110	6	8	4.8	3.2	3.2
120	9	6	4.8	1.2	2.2
80	5	3	4.8	-1.8	-2.8

These mistakes of model-1 will be given to model-2 as a 'target-variable'. Now model-2 will predict those mistakes based on given ZP 'gpa' & 'zq'.

→ In G.B. generally the second model is always Decision-Tree due to its performance in G.B.



→ This is how model-2 looks like.



$$\text{Prediction} = \text{Prediction-1} + \alpha * \text{Prediction-2}$$

actual - predicted

$$3 - (4.8 + 0.1 * (-1.8)) = -1.62$$

$$4 - (4.8 + 0.1 * (-0.8)) = -0.72$$

$$8 - (4.8 + 0.1 * (3.2)) = 2.88$$

$$6 - (4.8 + 0.1 * (2.8)) = 1.08$$

$$3 - (4.8 + 0.1 * (-2.8)) = -1.62$$

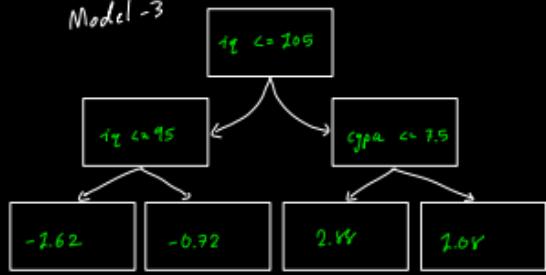
→ As we add models residual will move towards 0, α will keep moving.

- For Model-3 input will be $cypa$, iq but this model will predict the mistakes of first & second models.
- This process will continue for all n number of models, linked in this chain.



Prediction-3	Residual-3
(M3)	(M3)
-2.62	2.45
-0.92	-0.65
2.88	2.6
1.04	0.97
-2.62	2.46

Model -3

Residual-3 calculationactual - predicted

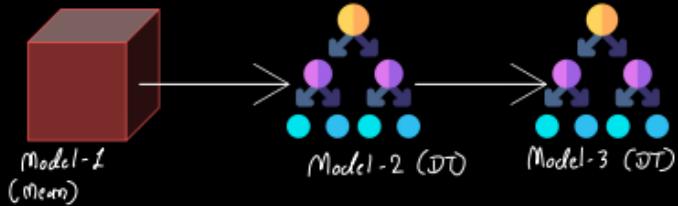
$$3 - (4.8 + 0.2 * (-2.8) + 0.2 * (-2.62)) = 2.45$$

$$4 - (4.8 + 0.2 * (-0.8) + 0.2 * (-0.92)) = -0.65$$

$$4 - (4.8 + 0.2 * (3.2) + 0.2 * (2.88)) = 2.6$$

$$6 - (4.8 + 0.2 * (2.2) + 0.2 * (1.04)) = 0.97$$

$$3 - (4.8 + 0.2 * (2.8) + 0.2 * (-2.62)) = -2.46$$

AdaBoost vs GradientBoost

1. max_leaf_node

[a]. AdaBoost

- ↳ We uses decision stumps
- ↳ Each tree will have maximum depth & f maximum leaf nodes = 2

[b]. GradientBoosting

- ↳ This is recommended that in regression task, max_leaf_node should be 8 to 32.

2. Learning_rate

[a]. AdaBoost

- ↳ Weight will be assigned to each model, based on their performance. (High performance = more weight)



[6]. GradientBoost

↳ Single constant weight will be multiplied with all model's prediction value.

Gradient Boosting Classifier

- Majority of process will be similar to GB regressor
- Process will slightly different to calculate initial prediction and class assignment.

	cgpa	iq	is_placed
0	6.82	118	0
1	6.36	125	1
2	5.39	99	1
3	5.50	106	1
4	6.39	148	0
5	9.13	148	1
6	7.17	147	1
7	7.72	72	0

→ for 1st model, the probability will be calculated using log-odd method.

$$\begin{aligned}\text{log(odd)} &= \log \left(\frac{\text{number of } 1\text{s in col column(Y)}}{\text{number of } 0\text{s in col column(Y)}} \right) \\ &= \log \left(\frac{5}{3} \right) \\ &= 0.51\end{aligned}$$

	cgpa	iq	is_placed	pre1(log-odds)
0	6.82	118	0	0.510826
1	6.36	125	1	0.510826
2	5.39	99	1	0.510826
3	5.50	106	1	0.510826
4	6.39	148	0	0.510826
5	9.13	148	1	0.510826
6	7.17	147	1	0.510826
7	7.72	72	0	0.510826

→ Now we need to calculate pseudo-residual (Actual - predicted.)

→ But we won't directly compute residual by subtracting pre1(log-odds) from is_placed because is_placed shows probability & pre1(log-odds) has log values.

→ Hence, we will convert pre1(log-odds) into probability using sigmoid.

$$\text{probability} = \left(\frac{1}{1 + e^{-\text{log(odd)}}} \right)$$

→ Now probability values are > 0.5 (assumed threshold) residual₁ = ($\text{is_placed} - \text{pre1(probability)}$)

→ In short, we have replaced mean calculation process (GBR for first model only)

with log(odd) for GB classification tasks.

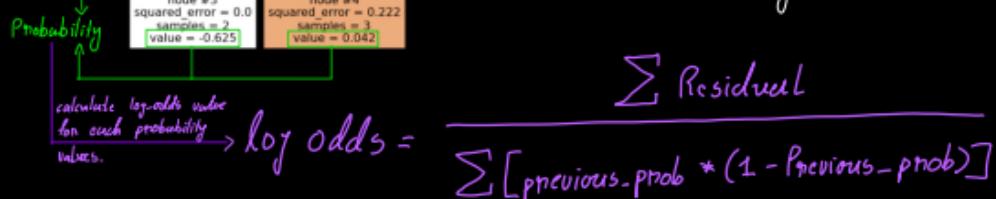


→ Now we will train a new DecisionTreeRegression model in which input will be `cpca`, `iq` and target variable will be `residual_1`.

→ This is how model-2 looks like.

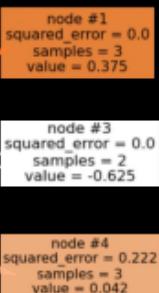
→ The leaf node of these decision tree returns probability value, hence we can't directly compare it with output of model-1 (log-odds value).

→ So `pre1(log-odds)` are in log format & leaf node of model-2 are in probability format so we need to convert these probability in log format. For all leaf node output values using below formula.



→ Now we will extract the Bucket address of each row.
(which row is passed through which decision leaf node)

	cpca	iq	is_placed	pre1(log-odds)	pre1(probability)	res1	leaf_entry1
0	6.82	118	0	0.510826	0.625	-0.625	3
1	6.36	125	1	0.510826	0.625	0.375	1
2	5.39	99	1	0.510826	0.625	0.375	1
3	5.50	106	1	0.510826	0.625	0.375	1
4	6.39	148	0	0.510826	0.625	-0.625	4
5	9.13	148	1	0.510826	0.625	0.375	4
6	7.17	147	1	0.510826	0.625	0.375	4
7	7.72	72	0	0.510826	0.625	-0.625	3



log odds for #node 3 (-0.625, -0.625)

$$= \frac{-0.625 - 0.625}{0.625(1-0.625) + 0.625(1-0.625)} = \frac{-1}{0.375}$$

$$= \frac{-2 \times 0.625}{0.625(1-0.625) + 0.625(1-0.625)} = \boxed{-2.66}$$

$$= \frac{-2 \times 0.625}{0.625(1-0.625) + 0.625(1-0.625)} = \frac{-2 \times 0.625}{0.625 \times 0.375 + 0.625 \times 0.375}$$

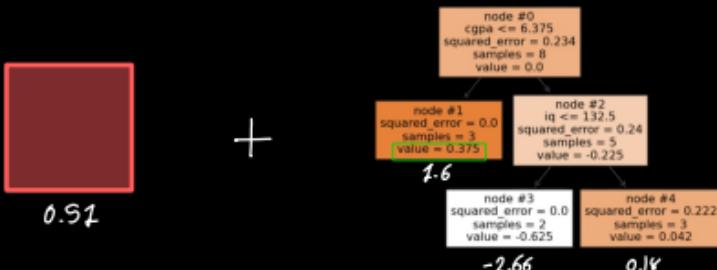
$$= \frac{-2 \times 0.625}{2 \times 0.625 \times 0.375}$$

$$\text{node } \#1 = 7.6$$

$$\text{node } \#3 = -2.66$$

$$\text{node } \#4 = 0$$

→ Now we will combine the o/p of both models.
→ Overall o/p log-odds for each row.



cgpa	iq	is_placed	pre1(log-odds)	pre1(probability)	res1	leaf_entry1	pre2(log-odds)
0	6.82	118	0	0.510826	0.625	-0.625	3
1	6.36	125	1	0.510826	0.625	0.375	1
2	5.39	99	1	0.510826	0.625	0.375	1
3	5.50	106	1	0.510826	0.625	0.375	1
4	6.39	148	0	0.510826	0.625	-0.625	4
5	9.13	148	1	0.510826	0.625	0.375	4
6	7.17	147	1	0.510826	0.625	0.375	4
7	7.72	72	0	0.510826	0.625	-0.625	3

cgpa	is_placed	is_placed	pre1(log-odds)	pre1(probability)	res1	leaf_entry1	pre2(log-odds)	pre1(probability)	res2	leaf_entry2
0	6.82	118	0	0.510826	0.625	-0.625	3	-2.159174	0.103477	-0.103477
1	6.36	125	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049
2	5.39	99	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049
3	5.50	106	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049
4	6.39	148	0	0.510826	0.625	-0.625	4	0.660026	0.666151	-0.666151
5	9.13	148	1	0.510826	0.625	0.375	4	0.660026	0.666151	0.333949
6	7.17	147	1	0.510826	0.625	0.375	4	0.659926	0.666151	0.333949
7	7.72	72	0	0.510826	0.625	-0.625	3	-2.159174	0.103477	-0.103477

convert back to probability

$$\frac{1}{1 + e^{-\text{log-odd}}}$$

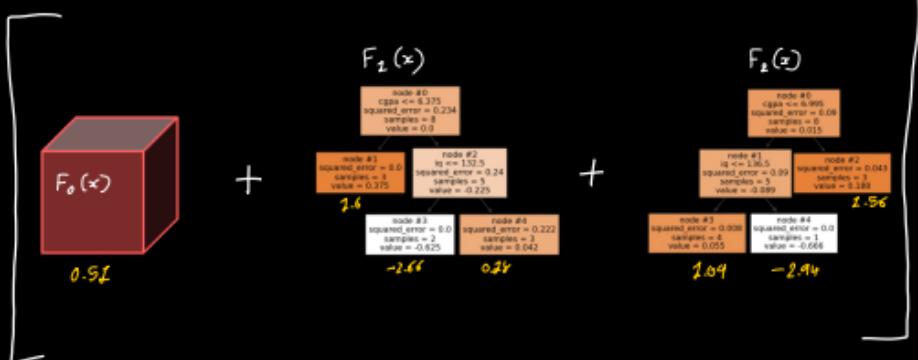
$$\hat{y} - \hat{y}^* = \uparrow$$

- Here we have noticed that the error is decreasing sharply it clearly indicates faster learning process.
- But we should not go with such amount of acceleration during learning.
- In order to reduce learning speed, we will use of (learning rate)



by multiplying with such log-odd values.

Now let's add one more model



	cgpa	iq	is_placed	pre1(log-odds)	pre1(probability)	res1	leaf_entry1	pre2(log-odds)	pre2(probability)	res2	leaf_entry2
0	6.82	118	0	0.510826	0.625	-0.625	3	-2.159174	0.103477	-0.103477	3
1	6.36	125	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049	3
2	5.39	99	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049	3
3	5.50	106	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049	3
4	6.39	148	0	0.510826	0.625	-0.625	4	0.690826	0.666151	-0.666151	4
5	9.13	148	1	0.510826	0.625	0.375	4	0.690826	0.666151	0.333849	2
6	7.17	147	1	0.510826	0.625	0.375	4	0.690826	0.666151	0.333849	2
7	7.72	72	0	0.510826	0.625	-0.625	3	-2.159174	0.103477	-0.103477	2

pre1(log-odds)	pre1(probability)
-1.068349	0.255717
3.201651	0.960896
3.201651	0.960896
3.201651	0.960896
-1.798349	0.142052
2.251651	0.904793
2.251651	0.904793
-0.698349	0.354722

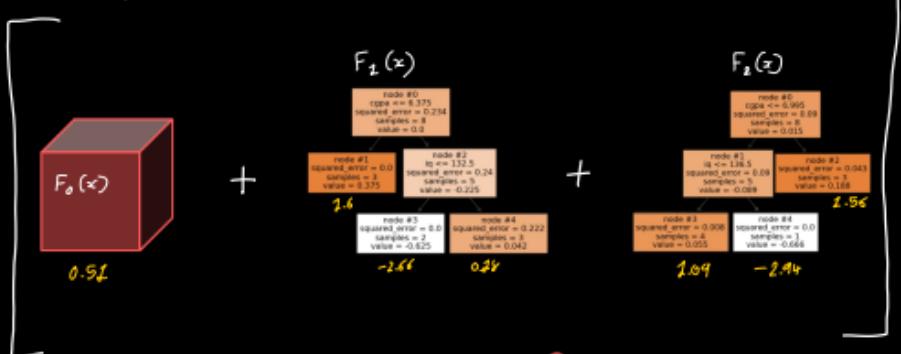
Now the training process is over, let's see how this model does prediction over unseen data.

* Steps

1. calculate log-odd values.
2. calculate probability of those log-odd values.
3. O/p based on probability.



Final trained model.

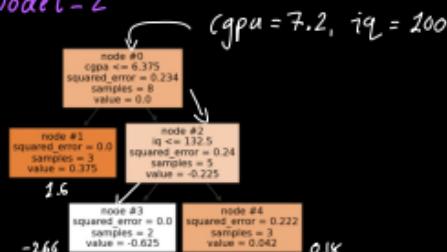


$$\text{query} = [[7.2, 100]] \rightarrow \{0, 1\} ??$$

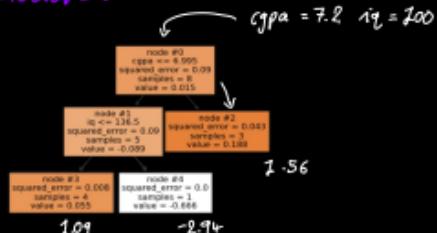
[a]. Model - 1

$$\text{log-odds} = 0.51$$

[b]. Model - 2



[c]. Model - 3



$$\text{log-odds} = 1.56$$

$$= 0.51 + (-2.66) + (1.56)$$

$$= \boxed{-0.59} \quad \text{log-odds}$$

$$\text{probability}(0.59) = \boxed{0.35}$$

assumed threshold = 0.5

$$\boxed{\text{olp} = 0}$$



XGBoost (Extreme Gradient Boosting)

- A ML model which can play very excellent on any kinds of dataset. No matter what is the length of characteristics of dataset.
- XGBoost is not any ml algorithm but it is a library - made over Gradient Boost with including some principals of software development.

Q: Why GradientBoost...?

1. flexible

↳ we can use any loss-functions.

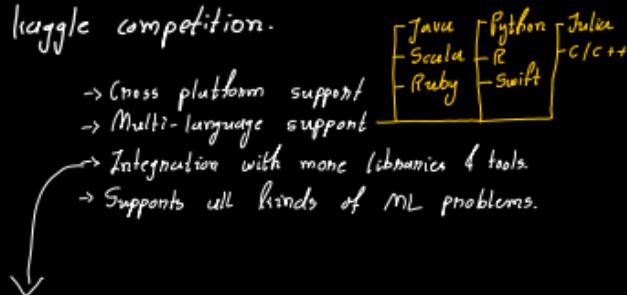
↳ we can use it for any kinds of problem like Regression, classification, Ranking, custom problems.

2. Performance

3. Robust

4. Most famous on Kaggle competition.

• Core areas of XGBoost



1. Performance

Compatibility

Python Libraries	Distributed Computing	Model interpretability	Model deployment	Workflow management
Numpy Pandas Matplotlib Scikit-learn	Spark PySpark Dask	SHAP Lime	Docker Kubernetes	Airflow mlflow

Usecases

- Regression, Classification, Time series forecasting, Ranking Problems, Anomaly Detection (Binary + multi-class)



Created with

Notewise

2. Speed

[a] Parallel processing

→ As we know in Gradient boosting, model building is sequential process.



→ The process will be in sequence, one model at a time.

→ As we know the process of building decision tree, where we try to find best splitting criterion.

Suppose we have a dataset with 2 zip & 7 colp variables like

type	is	placement
7.5	99	1
9.0	30	0
6.2	100	1

Now in order to get root node, we calculates best split of each column. after this process, we selects best one split from all these best splits of each columns.

• XGBoost uses parallel calculation mechanism, in which it calculates best split of each column simultaneously & then selects best one. hence the time consumption is less.

[b]. Optimized Data Structure

→ Generally ML models stores data row wise and processes data row by row.

→ XGBoost uses column block data structures in which data is stored in column format.

→ Each block will be a single column.

[c] cache awareness

→ Efficiently uses cache memory.

→ XGBoost uses histogram based binning.

→ XGBoost will create histogram for each column block and it stores range value of each bins in cache memory because those bin's range value will be used frequently by model.

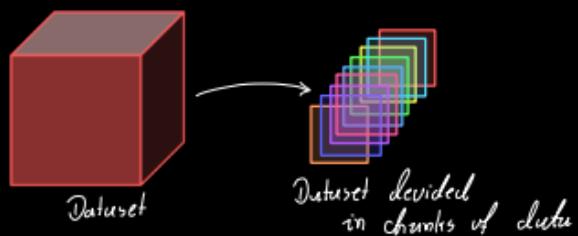
[d]. Out of core computing



↳ Data will be divided into chunks and then sequentially it loads each chunk of data.

Let's suppose we have a laptop with 8GB RAM, and you want to work with a dataset whose size is 45GB.

In this situation, XGBoost is here to help. by dividing dataset to small chunks.



[e]. Distributed Computing

↳ The process of dividing a single task between multiple machines.

↳ For XGBoost, we will distribute data between multiple machines & train the models parallelly and in last we will aggregate all models.

↳ For this, we need to use "DASK" or "Kubernetes"

[f]. GPU support

3. Performance

[a]. Regularizing learning objective

↳ In GB we does regularization using learning rate and pruning, but in XGBoost, it has in-build regularization feature in its loss function.

[b]. Sparsity Aware Split finding

↳ Sparsity : A dataset which contains too many zeros.

↳ XGBoost has capability to identify and handle sparsity.

Let's see how XGBoost does.

Sample dataset

f_2 gini-entropy

4 - 4.5

5 - 5.5

6 -

NaN - 7

8 -

9 - 4.5

best split

$f_2 > 4.5$

$f_2 > 4.5$

Yes

5, 6, 8, 9

No

4

→ Now based on above split, algorithm will put 'NaN' valued row in both leaf node one-by-one and calculate 'information-gain'. whichever node has higher 'information-gain', 'NaN' will be placed in that particular leaf node.

[c] Handling Missing Values

↳ Handles missing values internally

[d]. Efficient split finding (Weighted Quantile Sketch + Approximate Tree Learning)

f_2

3 - 4

5 - 6

7 - 8

9 - 10

11 - 12

13 - 12

} we calculates 'gini-entropy' using "exact-greedy-search"

method.

→ But as we can see, this approach is very

slow

→ To solve this problem Approximate Tree Learning approach comes into picture.

Approximate Tree Learning

↳ Convert this numerical column dataset into bins like [2-5], [6-10], [11-15]... (histogram based training)

→ Now this continuous variable became discrete variable.

→ This bins creation process is called as "histogram-based-training" for this process we uses a technique is called as Weighted Quantile Sketch

→ We selects best bin range based on 6th distribution

tion of the variable.

→ We should use "Quantile" approach to create bins ensuring each bin contains equal number of data points

[e]. Tree Pruning

→ The process of reducing the depth of a tree is called us pruning.

→ This method can helps to reduce the complexity of model. in short 'Tree pruning' reduces overfitting.

Types of pruning

→ Post-pruning

→ Pre-pruning

→ There are few more algorithms like LightGBM, CatBoost for better accuracy and robust performance.

