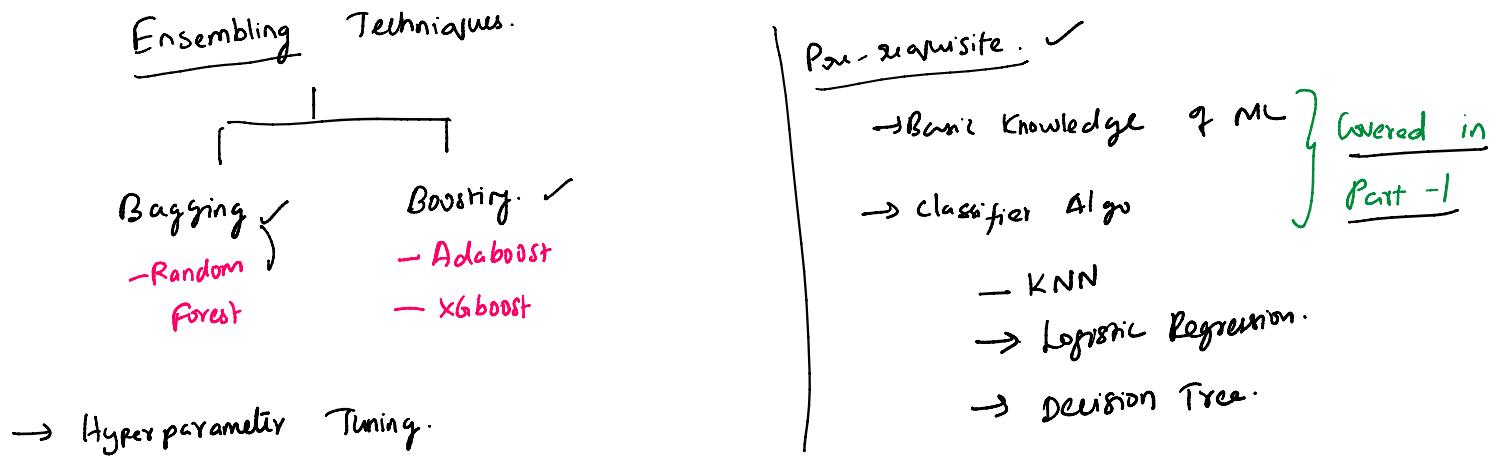


Introduction



Simple Ensembling. → Combining the result from diff. Models

Majority Voting. → Classifier.

- ① KNN → Non-spam
 - ② Logistic Regression → Spam
 - ③ Decision Tree. → Spam
- 2/3 → Spam.

Averaging. → Regression.

Advance techniques:

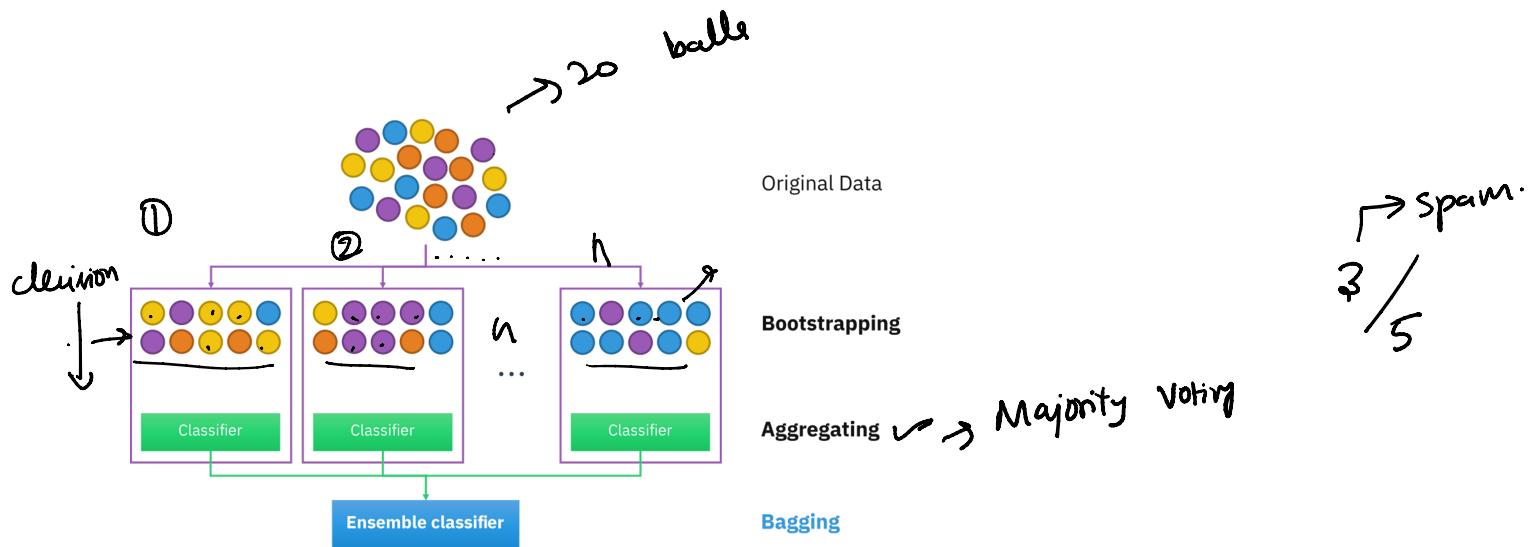
- ① Stacking.
- ② Bagging.
- ③ Boosting.

Bagging → Random Forest.

Bagging → Bootstrap + Aggregation.

Bootstrap

↳ Sampling with replacement.



Credit : https://commons.wikimedia.org/wiki/File:Ensemble_Bagging.svg

Random Forest

Eg: Titanic Dataset

original dataset → 4 Feature

Target

	Age	Gender	Class	Ticket	Survived
①					Y
②					Y
③					N
⋮					⋮
100					

→ Build decision tree.

Model ①

4 $\text{var}(y) = 2$

Age	Gender	Ticket	Survived

Source of Randomness:

- ① Bagging → data Sample.
- ② Features → \sqrt{n} / auto

How do Random forest Work?

- ① First built decision decision tree Model with each dataset → classifier.
- ② when a new test data point comes → feed the test data into each of the classifier
- ③ Use Majority Voting to finalize the prediction.

Model ②

	Gender	Class	Ticket	Survived
1				
2				
⋮				
5				

Model ③

	Age	Class	Ticket	Survived
1				
2				
⋮				
3				

Hyper parameters

n-estimator = No. of trees.

oob = out of bag sample.

oob → Test Score

Boosting

R.F

Bagging: i) The trees are independent.

ii) The result of one tree doesn't affect another tree. — Can run in parallel

Issue:

Bagging doesn't address the issues of previous tree

Base learner:

complete decision tree

Boosting:

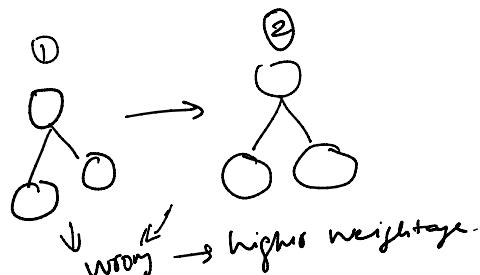
i) Boosting address the issue of bagging.

ii) It builds on the mistake of previous tree

iii) It runs sequentially and give more importance to mistakes
made by previous tree.

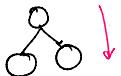
Base learner: stump / tree with 2-3 level deep

Weak learners.



Ada boost

→ Bare learner: Stump



→ Each stump has diff. say in the final decision.

→ At start each record has equal weightage.

→ Identify the best split [here is based on Weighted Gini impurity]

→ We check the prediction which went wrong.

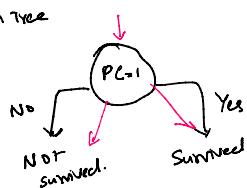
→ Change and give higher weightage to the data sample which are wrong.

→ Repeat the same process again with the new weight we calculated.

Lets go through an example:

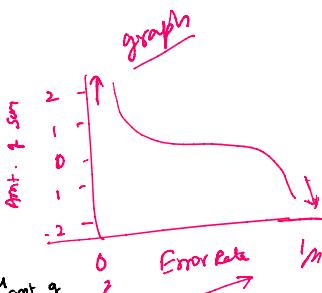
Passenger Class	Age	Gender	Survived	Initial Weight	Wt off 1st tree	Normalized weight
1	35	M	N	$\frac{1}{6}$	0.12	0.125
1	47	F	N	$\frac{1}{6}$	0.24	0.25
1	40	F	Y	$\frac{1}{6}$	0.12	0.125
3	45	M	N	$\frac{1}{6}$	0.12	0.125
3	35	F	N	$\frac{1}{6}$	0.12	0.125
3	55	M	Y	$\frac{1}{6}$	0.24	0.25
					0.96	n

① Decision Tree



② Error Rate = $\frac{\text{Wrong prediction}}{\text{Total Sample}} = \frac{2}{6}$

③ Amt of say = $\frac{1}{2} \log \left(\frac{1 - \text{ErrorRate}}{\text{ErrorRate}} \right)$



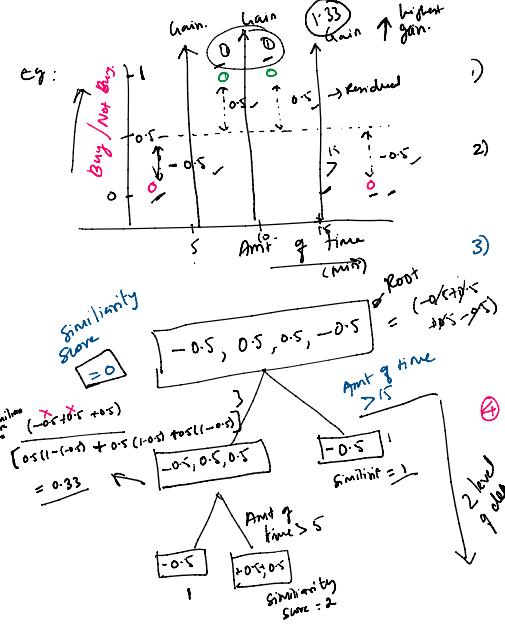
④ New weight

Incorrect Sample = Existing weight $\times e^{\frac{1}{2} \text{amt of say}}$

Correct Sample = Existing weight $\times e^{-\frac{1}{2} \text{amt of say}}$

5 Tree
 Survived (0.8)
 Non survive (0.9)
 ① 0.2
 ② 0.4
 ③ 0.2

The tree building process involved in XGBoost is different.

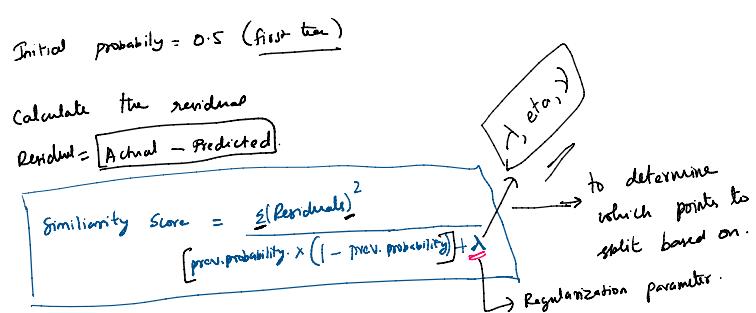
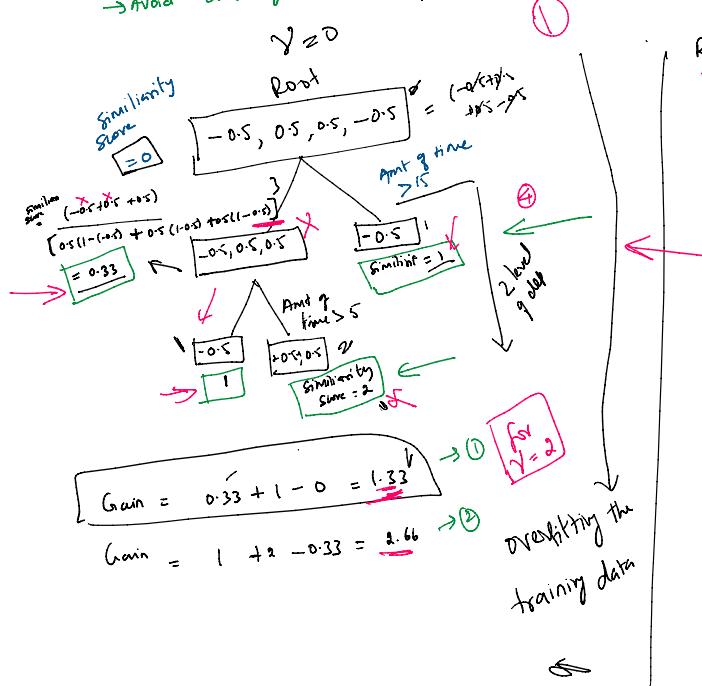


$$\text{Grain} = 0.33 + 1 - 0 = 1.33$$

$$\text{Grain} = 1 + 2 - 0.33 = 2.66$$

Gamma(γ) — pruning parameter
→ Avoid overfitting.

If $\gamma > \text{Grain}$
we prune the tree.



- Regularization parameter → λ
- ① Reduces similarity score / Grain.
 - ② Reduces overfitting → influence of single isolated value.

Regulation parameter → λ

Grain = $0.14 + 0.2 - 0 = 0.34$

Grain = $0.2 + 0.66 - 0.14 = 0.72$

$$1 \rightarrow 0.2 \quad (50\%) \quad \text{fitting part}$$

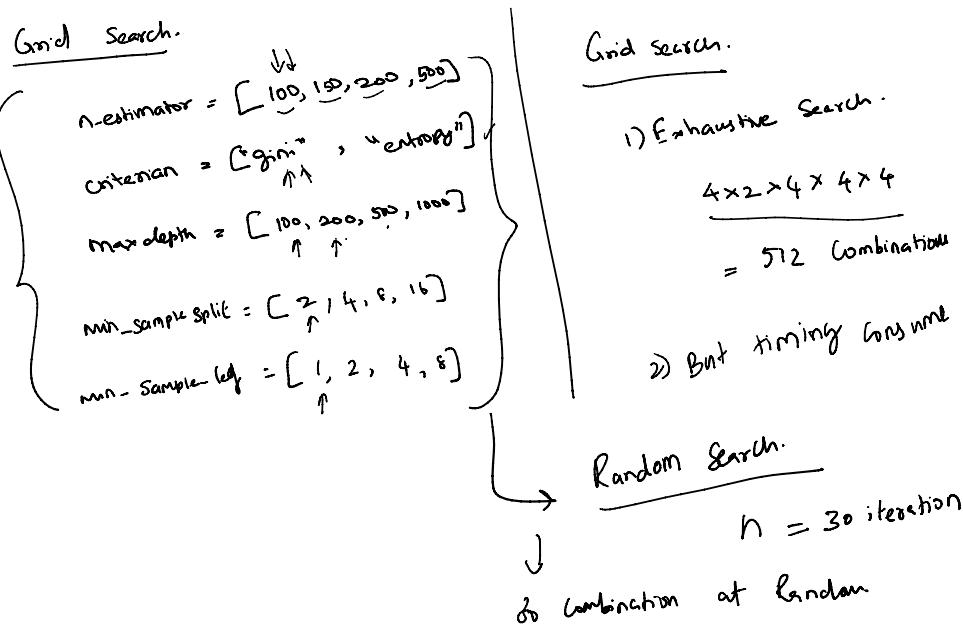
$$2 \rightarrow 0.66 \quad (66\%) \quad \text{fitting part}$$

Random, Grid Search & Bayesian Optimization

Parameters:

- n_estimators : int, default=100 → RF hyper-parameter.
The number of trees in the forest.
- Changed in version 0.22: The default value of n_estimators changed from 10 to 100 in 0.22.
- criterion : {"gini", "entropy"}, default="gini"
The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain. Note: this parameter is tree-specific.
- max_depth : int, None
The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.
- min_samples_split : int or float, default=2
The minimum number of samples required to split an internal node.
• If int, then consider min_samples_split as the minimum number.
• If float, then min_samples_split is a fraction and ceil(min_samples_split * n_samples) are the minimum number of samples for each split.
- Changed in version 1.8: Added float values for fractions.
- min_samples_leaf : int or float, default=1
The minimum number of samples required to be at a leaf node. A split point at any depth will only be considered if it leaves at least min_samples_leaf training samples in each of the left and right branches. This may have the effect of smoothing the model, especially in regression.
• If int, then consider min_samples_leaf as the minimum number.
• If float, then min_samples_leaf is a fraction and ceil(min_samples_leaf * n_samples) are the minimum number of samples for each node.

<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>



3) Bayesian Optimization.

1) Probabilistic model based approach.

→ Evidence based.

Objective function → "Loss function"

by minimize fn.
by changing the hyperparameter

2) start with Random hyper parameter.

↳ identify the space with good performance

↳ then choose other hyperparameter
based on the evidence found earlier.

Hyper Opt.

① parameter space

② fix the objective fn.

③ choose an Alp

→ Gaussian methods

→ Parzen Estimator etc...

Steps

① Start with default prediction $\boxed{0.5}$ \rightarrow first classifier.

② Calculate similarity score

$$\text{Similarity formula} = \left(\sum \text{residuals} \right)^2$$

Residual = Actual - predicted

$$\sum \text{prev probability} \times (1 - \text{prev probability}) + \Delta$$

③ Calculate Gain.

$$\text{Gain} = \frac{\text{Left similarity score}}{\text{Root similarity score}} + \frac{\text{Right similarity score}}{\text{Root similarity score}} \uparrow \text{MAX.}$$

④ Prune the tree if $\text{gamma}(\nu) > \text{Gain}$.

⑤ Calculate the output of the leaf using.

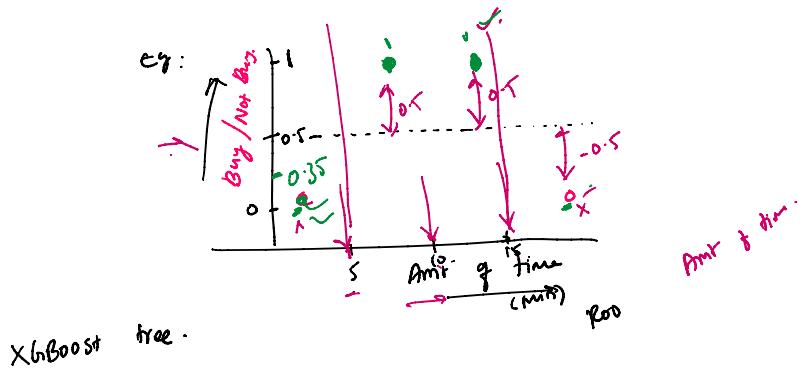
$$\boxed{\text{output} = \frac{\sum \text{Sum of Residuals}}{\sum \text{prev probabilities} \times (1 - \text{prev probability})}} \rightarrow \text{output in Log odds}$$

⑥ Final prediction

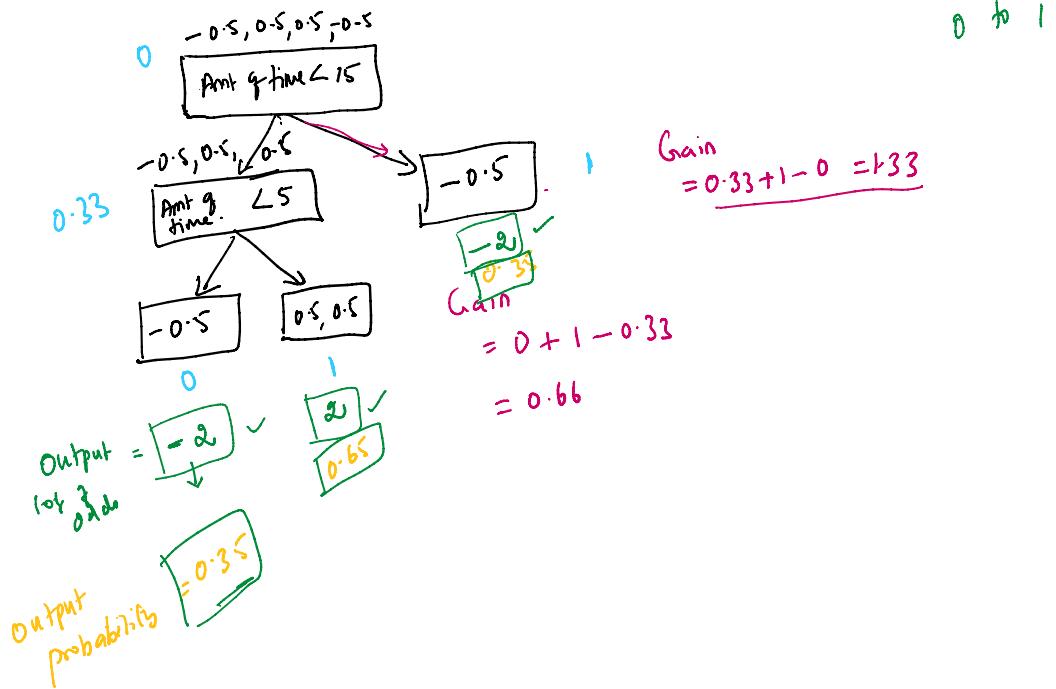
$$\begin{aligned} \text{output in log odds} &= \log \text{ odds} \text{ of previous probability} + \text{learning rate} (\text{eta}) \times \text{output of leaf} \\ &= 0 + 0.3 \times 2 = 0.6 \end{aligned}$$

$$\begin{aligned} \text{output probability} &= \frac{e^{\text{log odds}}}{1 + e^{\text{log odds}}} = \frac{e^{-0.6}}{1 + e^{-0.6}} = 0.35 \end{aligned}$$

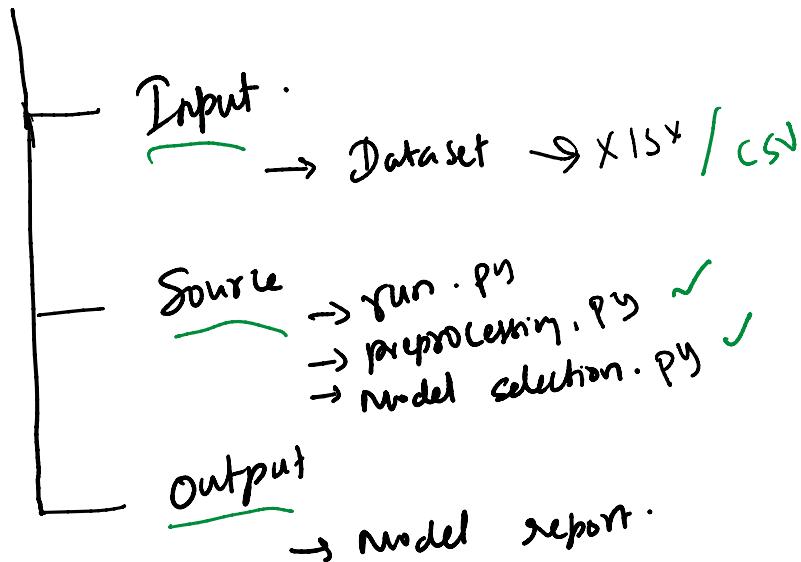
$$\begin{aligned} \log \left(\frac{P}{1-P} \right) &\rightarrow 0 \Rightarrow \log \left(\frac{0.35}{1-0.35} \right) = \log \left(\frac{0.35}{0.65} \right) = \log \left(\frac{1}{2} \right) = 0 \end{aligned}$$



XGBoost tree.



Conclusion & verdict on model performance



	model	precision_score	recall_score	f1_score	model_type
1	LogisticRegression	0.263339	0.242090	0.228505	Imbalanced
5	AdaBoostClassifier	0.542182	0.404654	0.437416	Imbalanced
2	GaussianNB	0.593930	0.749998	0.553084	Imbalanced
0	KNNClassifier	0.837556	0.833380	0.835141	Imbalanced
3	DecisionTreeClassifier	0.876315	0.920471	0.895335	Imbalanced
4	RandomForestClassifier	0.943305	0.944560	0.943913	Imbalanced
6	XGBoostClassifier	0.955169	0.942856	0.948765	Imbalanced