

Machine Learning Mastery

Worked Example

We can make the bootstrap procedure concrete with a small worked example. We will work through one iteration of the procedure.

Imagine we have a dataset with 6 observations:

```
1 [0.1, 0.2, 0.3, 0.4, 0.5, 0.6]
```

The first step is to choose the size of the sample. Here, we will use 4.

Next, we must randomly choose the first observation from the dataset. Let's choose 0.2.

```
1 sample = [0.2]
```

This observation is returned to the dataset and we repeat this step 3 more times.

```
1 sample = [0.2, 0.1, 0.2, 0.6]
```

We now have our data sample. The example purposefully demonstrates that the same value can appear zero, one or more times in the sample. Here the observation 0.2 appears twice.

An estimate can then be calculated on the drawn sample.

```
1 statistic = calculation([0.2, 0.1, 0.2, 0.6])
```

Those observations not chosen for the sample may be used as out of sample observations.

```
1 oob = [0.3, 0.4, 0.5]
```

In the case of evaluating a machine learning model, the model is fit on the drawn sample and evaluated on the out-of-bag sample.

```
1 train = [0.2, 0.1, 0.2, 0.6]
2 test = [0.3, 0.4, 0.5]
3 model = fit(train)
4 statistic = evaluate(model, test)
```

That concludes one repeat of the procedure. It can be repeated 30 or more times to give a sample of calculated statistics.

```
1 statistics = [...]
```

This sample of statistics can then be summarized by calculating a mean, standard deviation, or other summary values to give a final usable estimate of the statistic.

```
1 estimate = mean([...])
```

1. Choose a number of bootstrap samples to perform
2. Choose a sample size
3. For each bootstrap sample
 1. Draw a sample with replacement with the chosen size
 2. Fit a model on the data sample
 3. Estimate the skill of the model on the out-of-bag sample.
4. Calculate the mean of the sample of model skill estimates.

```
from sklearn.utils import resample
```

```
data = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6]
data
```

```
[0.1, 0.2, 0.3, 0.4, 0.5, 0.6]
```

```
## sample observations
sample = resample(data, replace=True, n_samples=4, random_state=1)
sample
```

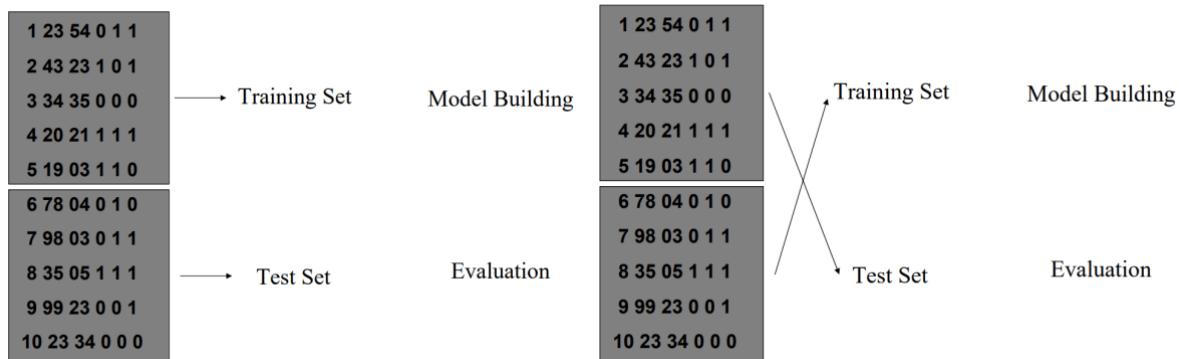
```
[0.6, 0.4, 0.5, 0.1]
```

```
# out of bag observations
oob = [x for x in data if x not in sample]
oob
```

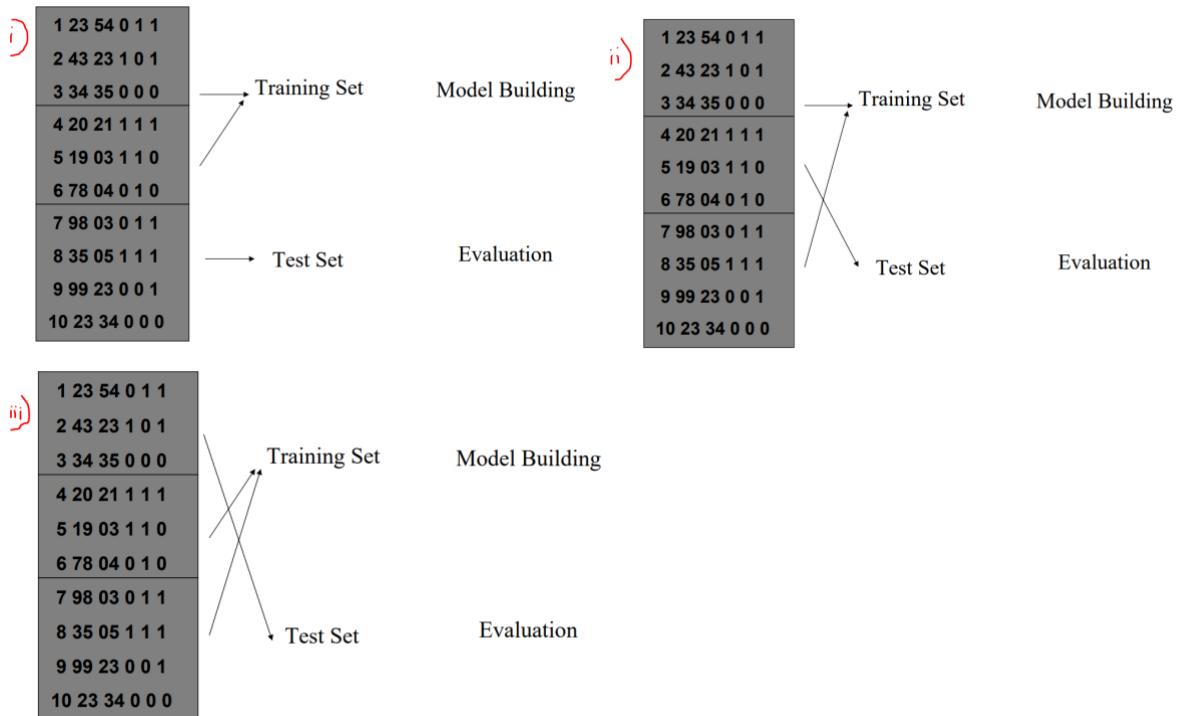
```
[0.2, 0.3]
```

Cross-Validation

Leave N/2 out



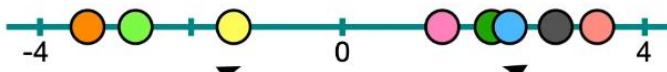
Leave N/3 out



Stats-Quest

Reason for Bootstrapping

Feeling Worse ← → Feeling Better



...and we gave that drug to 8 different people that had the illness.



vs

There is a new drug to treat illness.



Feeling Worse ← → Feeling Better

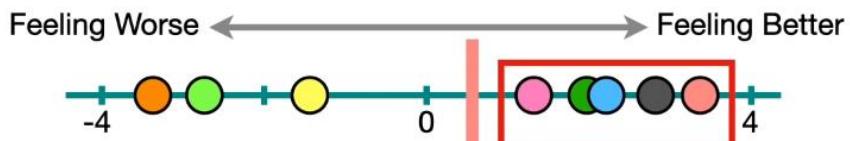


For 5 of those people, the drug appeared to help them feel better...

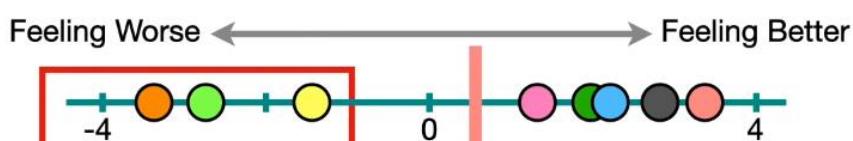
Feeling Worse ← → Feeling Better



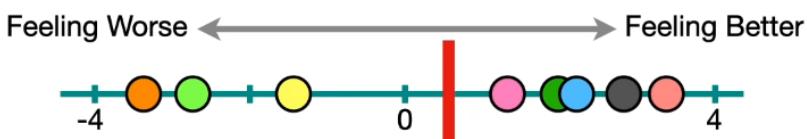
...but for 3 people, the drug appeared to make them feel worse.



However, maybe these **5** people all felt better because they were healthier to begin with...



...and maybe these **3** people all felt worse because they had unhealthy lifestyles.



*) So it is possible that the reason we got a mean value = **0.5** instead of **0** is because of random things that we can't control.

*) Is that can we do to decide the drug will work / not ?

The solution is

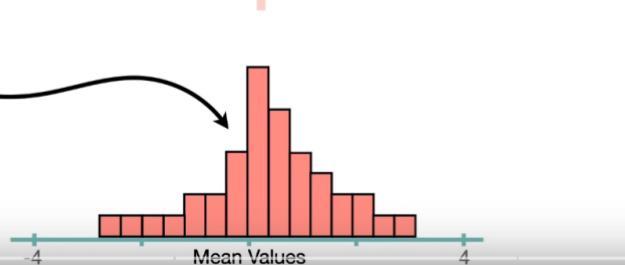
One **expensive** and **time consuming** option would be to replicate the experiment a bunch of times.



Feeling Worse ← → Feeling Better

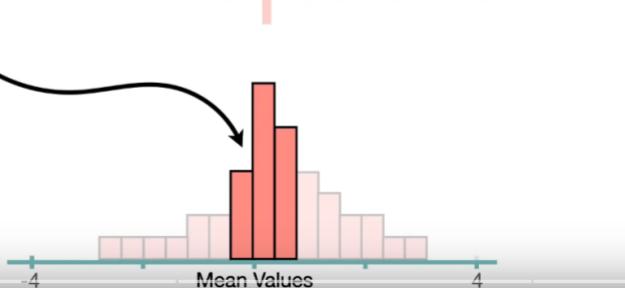
Repeating the experiments
a bunch of time and keep
track of the mean

...and we will end
up with a histogram
of mean values.



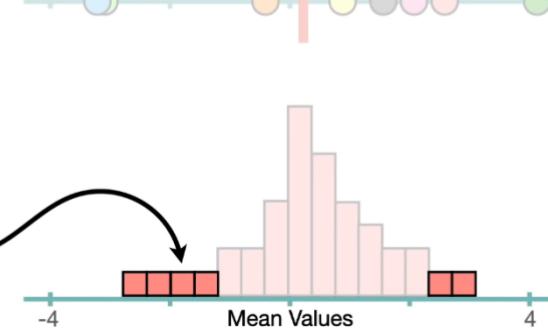
Feeling Worse ← → Feeling Better

Just by looking at
this distribution, we
can see that mean
values close to **0**,
which suggest the
drug does not do
anything, are
relatively likely to
occur...



Feeling Worse ← → Feeling Better

...and mean values
far from **0**,
indicating that the
drug does
something, are
relatively rare.



Repeating the experiments a bunch of time is expensive and time consuming,
one optimal solution to solve this is to use Bootstrapping.

Bootstrap in action

We regularly come across many game shows on television and you must have noticed an option of “Audience Poll”. Most of the time a contestant goes with the option which has the highest vote from the audience and most of the time they win. We can generalize this in real life as well where taking opinions from a majority of people is much more preferred than the opinion of a single person. The Ensemble technique has a similar underlying idea where we aggregate predictions from a group of predictors, which may be classifiers or regressors, and most of the time the prediction is better than the one obtained using a single predictor.

Definition: — Ensemble learning is a machine learning paradigm where multiple models (often called “weak learners”) are trained to solve the same problem and combined to get better results. The main hypothesis is that when weak models are correctly combined, we can obtain more accurate and/or robust models.

Weak Learners: A ‘weak learner’ is any ML algorithm (for regression/classification) that provides an accuracy slightly better than random guessing.

In ensemble learning theory, we call weak learners (or base models) models that can be used as building blocks for designing more complex models by combining several of them. Most of the time, these basic models perform not so well by themselves either because they have a high bias or because they have too much variance to be robust. Then, the idea of ensemble methods is to try reducing bias and/or variance of such weak learners by combining several of them together to create a strong learner (or ensemble model) that achieves better performances.

Let's suppose we have 'n' predictors/models:

Z₁, Z₂, Z₃,, Z_n with a standard deviation of σ

$$\text{Variance}(z) = \sigma^2$$

If we use single predictors Z₁, Z₂, Z₃,, Z_n the variance associated with each will be σ^2 but the expected value will be the average of all the predictors.

Let's consider the average of the predictors:

$$\mu = (Z_1 + Z_2 + Z_3 + \dots + Z_n)/n$$

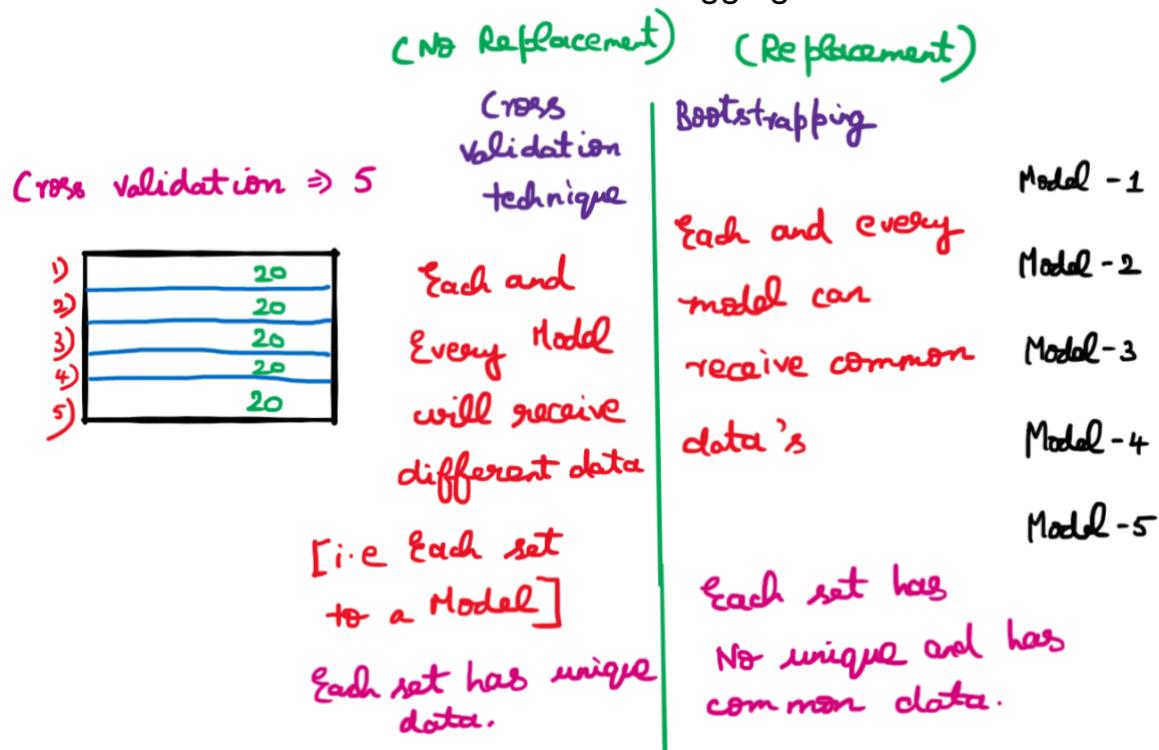
If we use μ as the predictor then the expected value remains the same but see the variance now:

$$\text{variance}(\mu) = \sigma^2/n$$

So, the expected value remained ' μ ' but variance decreases when we use an average of all the predictors.

This is why taking mean is preferred over using single predictors.

Difference between cross-validation and bagging



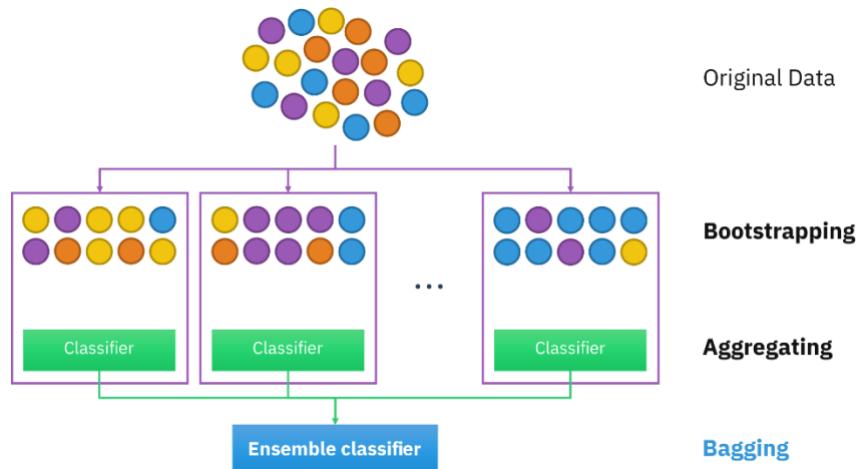
- Brief answer

Bootstrap aggregating (also called Bagging) is a method used to produce a predictive model. Cross-validation is a model validation technique. So cross-validation and bagging solve different problems. You cannot replace cross-validation with bagging as it would not have any sense.

- Explanation

Bagging consists in sampling several datasets from an original training set. Then, a model (generally a classification tree or a regression tree) is trained on each new dataset. Next, the predictions of all the models are aggregated (for example, by taking their mean for a regression task or the majority vote for a classification task) to get a final model.

Below is a schema of how bagging works.

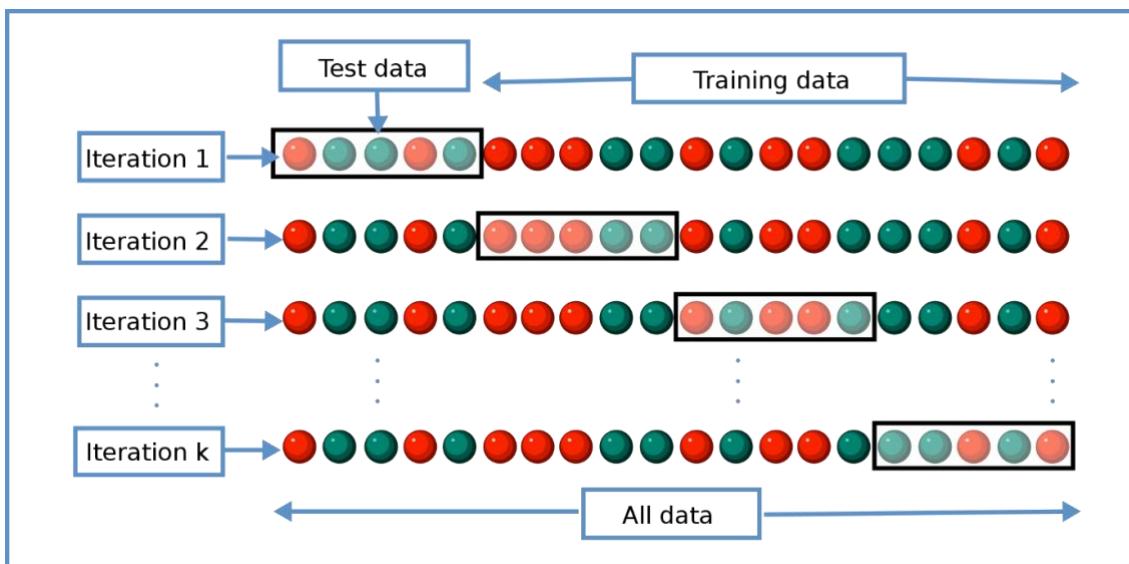


The main advantage of using bagging is indeed to reduce the variance. A model obtained by aggregating several models trained on several sampled dataset is generally less prone to overfitting.

On another side, cross-validation is a method to measure the reliability of a model once it is already created. For example, let's describe how works 5-fold cross-validation.

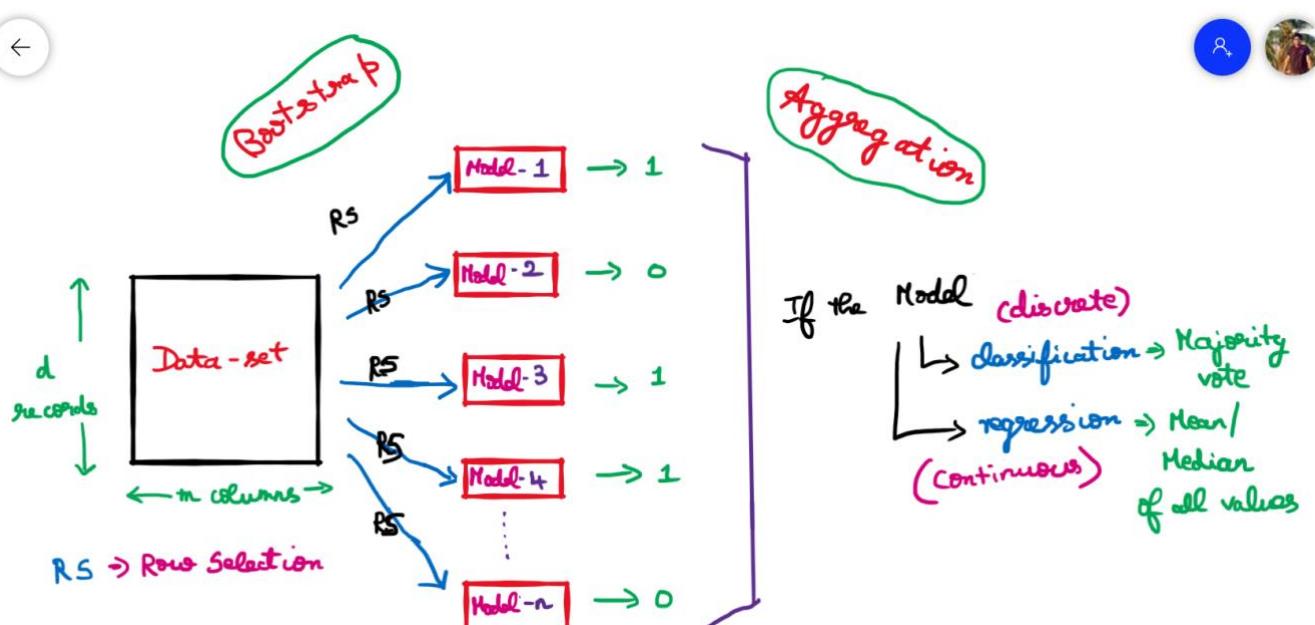
First, the original dataset is divided into 5 subsets. Next, the model is trained on four of them and a validation performance is computed on the fifth one. Then the operation is repeated by selecting another validation subset among the predefined subsets. At the end of the procedure, we have 5 performance scores, one per validation subset. The mean and standard deviation of the 5 performance scores can be calculated to estimate the bias and variance of the model.

Below is a schema of how k-cross-validation works (in our previous example, $k = 5$).



Bagging / Bootstrap Aggregation

It is a parallel process.



1. BAGGING

Bootstrap Aggregation

* Bagging stands for Bootstrap Aggregation.

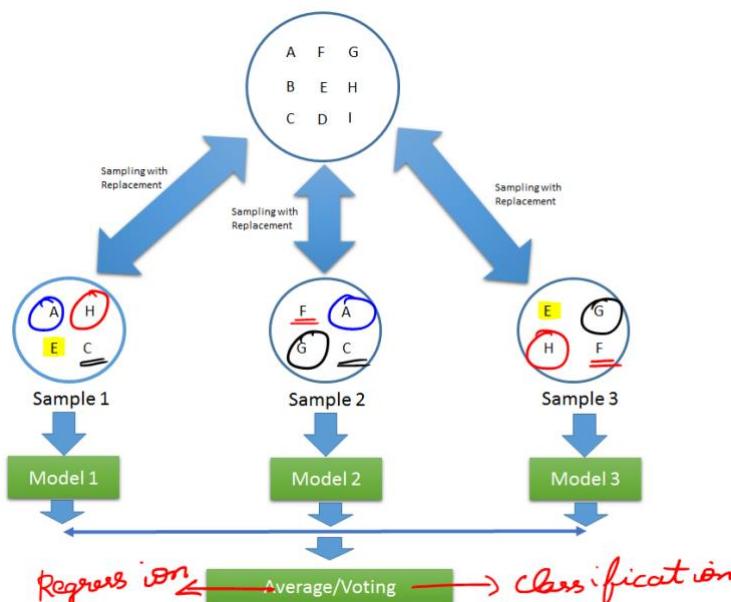
* In real-life scenarios, we don't have multiple different training sets on which we can train our model separately and at the end combine their result. Here, bootstrapping comes into the picture.

* Bootstrapping is a technique of sampling different sets of data from a given training set by using replacement. After bootstrapping the training dataset, we train the model on all the different sets and aggregate the result. This technique is known as Bootstrap Aggregation or Bagging.

* For aggregating the outputs of base learners, bagging uses majority voting (most frequent prediction among all predictions) for classification and averaging (mean of all the predictions) for regression.

Regression \Rightarrow Averaging

Classification \Rightarrow Majority



(with replacement)

· Advantages of a Bagging Model:

Since the model sees
repeated data
↑

1. Bagging significantly decreases the variance without increasing bias.
2. Bagging methods work so well because of diversity in the training data since the sampling is done by bootstrapping.
3. Also, if the training set is very huge, it can save computational time by training the model on a relatively smaller data set and still can increase the accuracy of the model.
4. Works well with small datasets as well.

· Disadvantages of a Bagging Model:

1. The main disadvantage of Bagging is that it improves the accuracy of the model at the expense of interpretability i.e., if a single tree was being used as the base model, then it would have a more attractive and easily interpretable diagram, but with the use of bagging this interpretability gets lost.
2. Another disadvantage of Bootstrap Aggregation is that during sampling, we cannot interpret which features are being selected i.e., there are chances that some features are never used, which may result in a loss of important information.

Out of Bag Evaluation: In bagging, when different samples are collected, no sample contains all the data but a fraction of the original dataset. There might be some data that are never sampled at all. The remaining data which are not sampled are called out of bag instances.

The Random Forest approach is a bagging method where deep trees (Decision Trees), fitted on bootstrap samples, are combined to produce an output with lower variance.

Pasting

Pasting is an ensemble technique similar to bagging with the only difference being that there is no replacement done while sampling the training dataset. This causes less diversity in the sampled datasets and data ends up being correlated. That's why bagging is more preferred than pasting in real scenarios.

Let's see python implementation of Bagging:

```
from sklearn.ensemble import BaggingClassifier  
from sklearn.neighbors import KNeighborsClassifier
```

```
from sklearn.datasets import load_breast_cancer  
dataset = load_breast_cancer()  
X = dataset.data  
y = dataset.target
```

```
from sklearn.model_selection import train_test_split  
X_train, X_test, y_train, y_test = train_test_split(  
    X, y, random_state=3  
)
```

```
knn = KNeighborsClassifier(n_neighbors=5)  
knn.fit(X_train, y_train)  
knn.score(X_test, y_test)
```

0.916083916083916

Bagging/bootstrap aggregation

```
## number of classifiers → 10  
## % of data each classifier can fit → 50%  
## bootstrap=False (i.e pasting)  
  
bag_knn = BaggingClassifier(KNeighborsClassifier(n_neighbors=5),  
                           n_estimators=10, max_samples=0.5,  
                           bootstrap=True, random_state=3, oob_score=True)
```

```
bag_knn.fit(X_train, y_train)  
bag_knn.score(X_test, y_test)
```

0.9370629370629371

Great! our score significantly improves with use of bagging.

let's not use bootstrap and see the model accuracy! Remember this is "Pasting"

Pasting

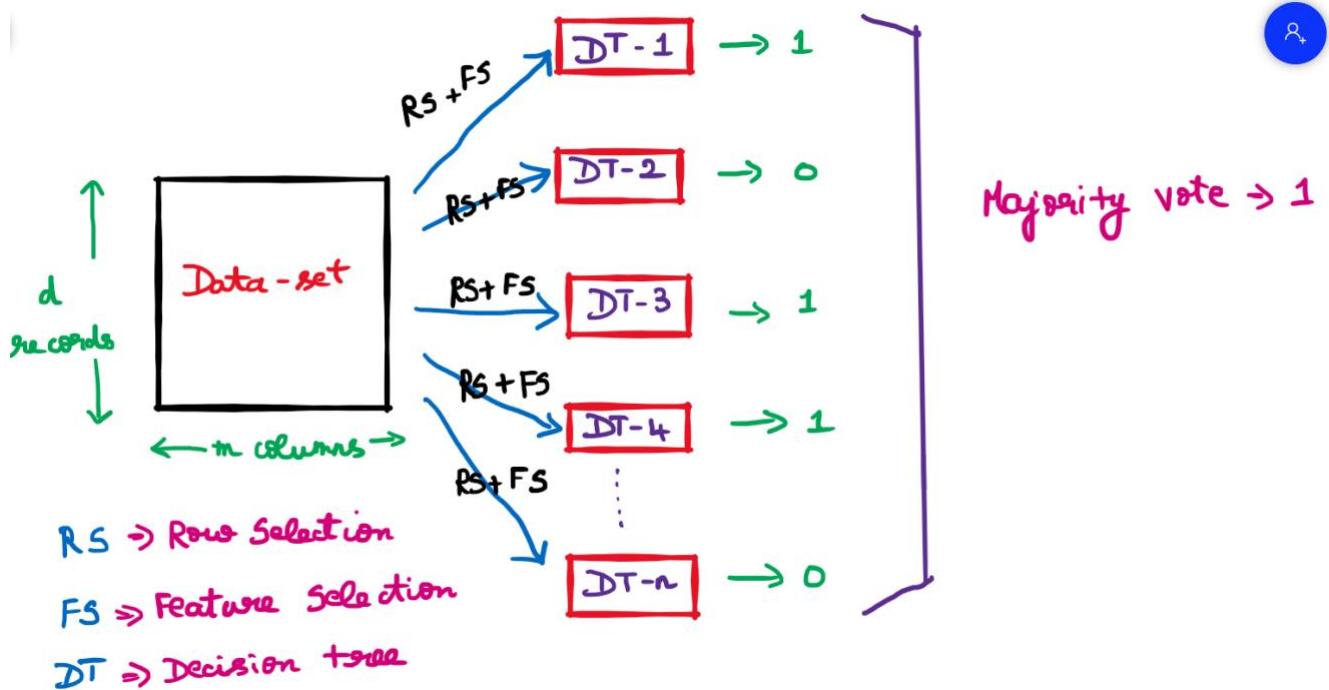
```
pasting_knn = BaggingClassifier(KNeighborsClassifier(n_neighbors=5),  
                                n_estimators=10, max_samples=0.5,  
                                bootstrap=False, random_state=3)
```

```
pasting_knn.fit(X_train, y_train)  
pasting_knn.score(X_test, y_test)
```

0.9300699300699301

Random Forest

The base learner is decision tree.



Advantages and Disadvantages of Random Forest:

- 1) It can be used for both regression and classification problems.
- 2) Since base model is a tree, handling of missing values is easy.
- 3) It gives very accurate result with very low variance.
- 4) Results of a random forest are very hard to interpret in comparison with decision trees.
- 5) High computational time than other respective models.

Random Forest should be used where accuracy is up utmost priority and interpretability is not very important. Also, computational time is less expensive than the desired outcome.

```
data = pd.read_csv("winequality_red.csv")
data.head()
```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
0	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5
1	7.8	0.88	0.00	2.6	0.098	25.0	67.0	0.9968	3.20	0.68	9.8	5
2	7.8	0.76	0.04	2.3	0.092	15.0	54.0	0.9970	3.26	0.65	9.8	5
3	11.2	0.28	0.56	1.9	0.075	17.0	60.0	0.9980	3.16	0.58	9.8	6
4	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5

The data set consists following Input variables :

1 - fixed acidity 2 - volatile acidity 3 - citric acid 4 - residual sugar 5 - chlorides
6 - free sulfur dioxide 7 - total sulfur dioxide 8 - density 9 - pH 10 - sulphates 11 - alcohol.

Output variable gives the quality of th wine based on the input variables:

12 - quality (score between 0 and 10)

```
X = data.drop(columns = 'quality')
y = data['quality']

x_train,x_test,y_train,y_test = train_test_split(X,y,test_size = 0.30, random_state= 355)
```

```
#let's first visualize the tree on the data without doing any pre processing
clf = DecisionTreeClassifier(criterion='gini',min_samples_split= 2)
clf.fit(x_train,y_train)
```

DecisionTreeClassifier()

```
# accuracy of our classification-1 tree
clf.score(x_test,y_test)
```

0.6291666666666667

```
#let's first visualize the tree on the data without doing any pre processing
clf2 = DecisionTreeClassifier(criterion = 'entropy', max_depth =24, min_samples_leaf= 1)
clf2.fit(x_train,y_train)
```

DecisionTreeClassifier(criterion='entropy', max_depth=24)

```
# accuracy of our classification-2 tree
clf2.score(x_test,y_test)
```

0.6208333333333333

```
rand_clf = RandomForestClassifier(criterion='gini',n_estimators=100,random_state=6)
```

Random state, if given none then score will vary everytime you run the RandomForestClassifier. If we asssign a value to it, then result will remain constant.

```
rand_clf.fit(x_train,y_train)

RandomForestClassifier(random_state=6)

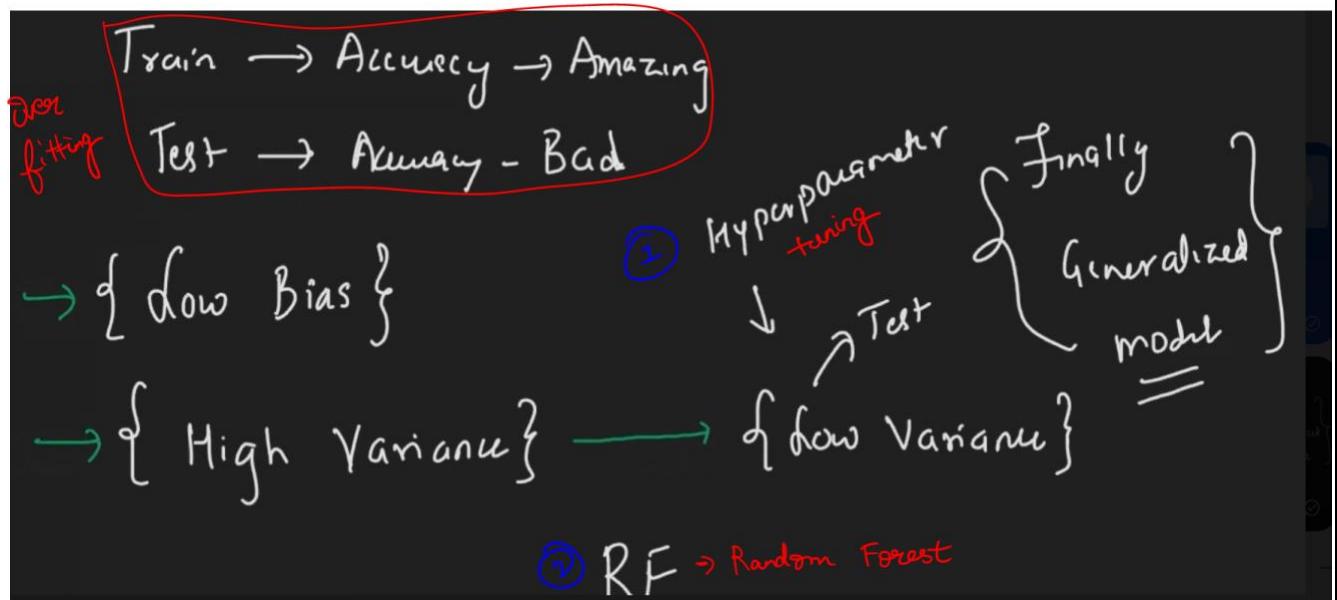
rand_clf.score(x_test,y_test)

0.6708333333333333
```

Is Random Forest impacted by the outliers ??? No

Is Random Forest requires feature scaling ?? No

Is over-fitting takes place in Random Forest ?? No



Hold-out Method

Name is different
Normal train-test-split method

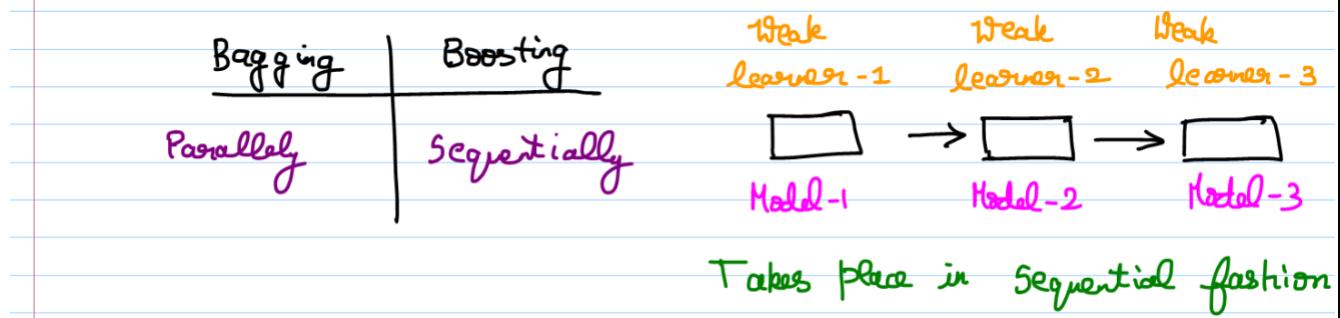
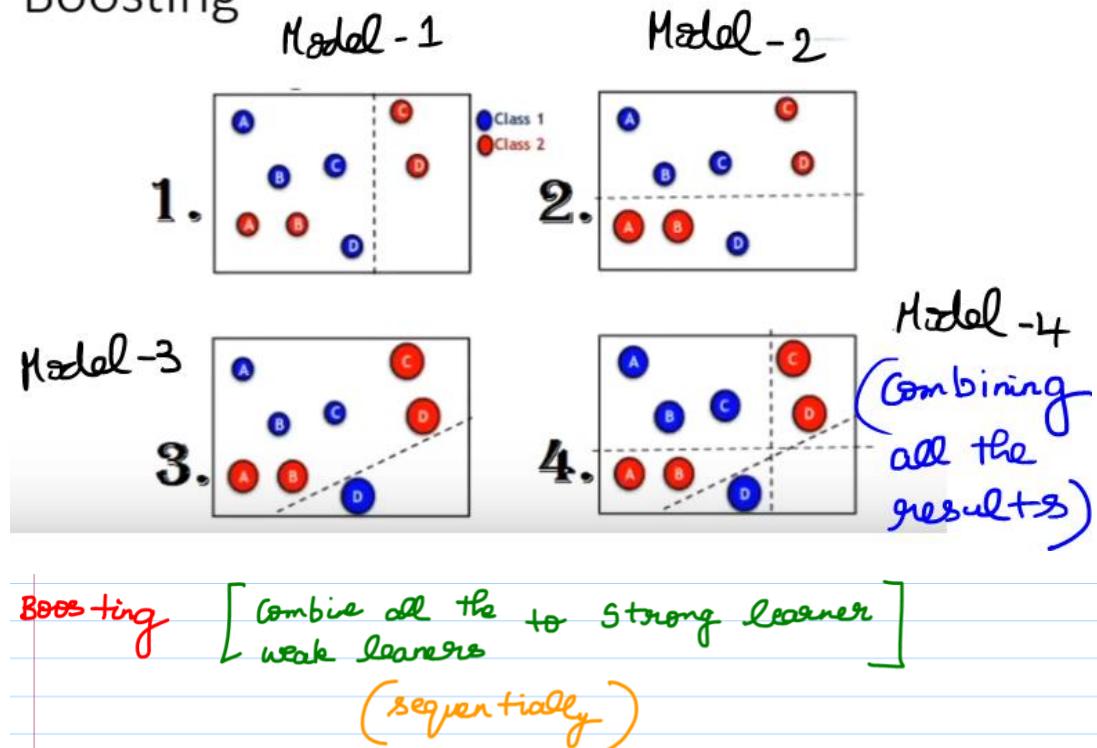
- Hold Out Method:

It is the most basic of the CV techniques. It simply divides the dataset into two sets of training and test. The training dataset is used to train the model and then test data is fitted in the trained model to make predictions. We check the accuracy and assess our model on that basis. This method is used as it is computationally less costly. But the evaluation based on the Hold-out set can have a high variance because it depends heavily on which data points end up in the training set and which in test data. The evaluation will be different every time this division changes.

Boosting

It is a sequential process.

Boosting



Ada-boost (Adaptive Boosting)

Steps for AdaBoost Algorithm

Step 1 : initialize the weights as 1/n to every n observations

Step 2 : Select the 1 features according to Lowest Gini/Highest information Gain and calculate the Total Error

Step 3 : calculate the Performance of the stump.

Step 4: Calculate the new weights for each misclassification (increase) and right classification (decrease)

Step 5 : Normalize the new weights so that the sum of weights is 1

Step 6: Now Repeat from Step 2 and so on till the configured number of estimators reached or the accuracy achieved.

Step-1 → Pick out the 1st decision stump

Learning An AdaBoost Model From Data

AdaBoost is best used to boost the performance of decision trees on binary classification problems.

AdaBoost was originally called AdaBoost.M1 by the authors of the technique Freund and Schapire. More recently it may be referred to as discrete AdaBoost because it is used for classification rather than regression.

AdaBoost can be used to boost the performance of any machine learning algorithm. It is best used with weak learners. These are models that achieve accuracy just above random chance on a classification problem.

The most suited and therefore most common algorithm used with AdaBoost are decision trees with one level. Because these trees are so short and only contain one decision for classification, they are often called decision stumps.

Each instance in the training dataset is weighted. The initial weight is set to:

$$\text{weight}(x_i) = 1/n$$

Where x_i is the i'th training instance and n is the number of training instances.

Data-set

	f_1	f_2	f_3	output	Initial weights
1)				Yes	$\frac{1}{7}$
2)				No	$\frac{1}{7}$
3)				Yes	$\frac{1}{7}$
4)				Yes	$\frac{1}{7}$
5)				No	$\frac{1}{7}$
6)				No	$\frac{1}{7}$
7)				Yes	$\frac{1}{7}$

Decision tree is created only with the help of 1-depth } Stumps
 with many leaf nodes }

*) Based on f_1



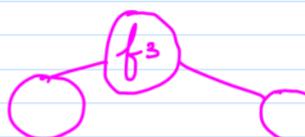
\Rightarrow Entropy / Gini Impurity }

*) Based on f_2



\Rightarrow Entropy / Gini Impurity }

*) Based on f_3



\Rightarrow Entropy / Gini Impurity }

} Take the lowest one

*) Now we selected 1st Decision tree as the 1st base learning model.
 (lowest entropy)

Step-2 → Calculate the error

A weak classifier (decision stump) is prepared on the training data using the weighted samples. Only binary (two-class) classification problems are supported, so each decision stump makes one decision on one input variable and outputs a +1.0 or -1.0 value for the first or second class value.

The misclassification rate is calculated for the trained model. Traditionally, this is calculated as:

$$\text{error} = (\text{correct} - N) / N$$

Where error is the misclassification rate, correct are the number of training instance predicted correctly by the model and N is the total number of training instances. For example, if the model predicted 78 of 100 training instances correctly the error or misclassification rate would be $(78-100)/100$ or 0.22.

Data-set

	f_1	f_2	f_3	output	Initial weights
1)				Yes	$\frac{1}{7}$
2)				No	$\frac{1}{7}$
3)				Yes	$\frac{1}{7}$
4)				Yes	$\frac{1}{7}$
5)				No	$\frac{1}{7}$
6)				No	$\frac{1}{7}$
7)				Yes	$\frac{1}{7}$

incorrectly classified by the 1st decision stump.

(Based on feature - 1)

$$\text{Total error} \Rightarrow \frac{\text{incorrect}}{N}$$

correct \Rightarrow correctly classified samples

N \Rightarrow Total Number of samples

$$\begin{array}{c|c} \text{correct} & \text{Incorrect} \\ \hline 6 & 1 \end{array} \Rightarrow \frac{6-1}{7} \Rightarrow \frac{1}{7}$$

$$\boxed{\text{Total error} \Rightarrow \frac{1}{7}}$$

$$\left. \begin{array}{l} \text{Performance of the stump} \\ (\text{stage}) \end{array} \right\} \Rightarrow \frac{1}{2} \log_e \left(\frac{1 - \text{Total error}}{\text{Total error}} \right)$$

$$\Rightarrow \frac{1}{2} \log_e \left(\frac{1 - \frac{1}{7}}{\frac{1}{7}} \right)$$

$$\Rightarrow \frac{1}{2} \log_e \left(\frac{\frac{6}{7}}{\frac{1}{7}} \right)$$

$$\Rightarrow \frac{1}{2} \log_e (6)$$

$$\boxed{\text{Performance of the stump} \Rightarrow 0.896}$$

Weight will be increasing for incorrectly classified sample.

Weight will be decreasing for correctly classified sample.

Update the weights

Correctly classified sample

New sample weight } \Rightarrow Old sample weight + $e^{-\text{Performance of the stump}}$

Incorrectly classified sample

New sample weight } \Rightarrow Old sample weight + $e^{\text{Performance of the stump}}$

Correctly classified sample } $\Rightarrow \frac{1}{7} * e^{-0.895} \Rightarrow 0.05$

Incorrectly classified sample } $\Rightarrow \frac{1}{7} * e^{0.895} \Rightarrow 0.349$

Step-3 Getting the Normalized weights

Data-set

	f_1	f_2	f_3	output	Initial weights	Updated weights	Normalized weights
1)				Yes	$\frac{1}{7}$	0.05	0.07
2)				No	$\frac{1}{7}$	0.05	0.07
3)				Yes	$\frac{1}{7}$	0.349	0.513
4)				Yes	$\frac{1}{7}$	0.05	0.07
5)				No	$\frac{1}{7}$	0.05	0.07
6)				No	$\frac{1}{7}$	0.05	0.07
7)				Yes	$\frac{1}{7}$	0.05	0.07

Since the sum of weights!=1, we are normalizing the weights.

Normalized weights \Rightarrow Updated weights / 0.68

$$0.05 / 0.68 \Rightarrow 0.07$$

$$0.349 / 0.68 \Rightarrow 0.513$$

Step-4 Creating the New data-set

	Data-set		Normalized weights	Making buckets		
	f_1	f_2	f_3	output		
1)				Yes	0.07	0 to 0.07
2)				No	0.07	0.07 to 0.14
3)				Yes	0.513	0.14 to 0.653
4)				Yes	0.07	0.653 to 0.723
5)				No	0.07	0.723 to 0.793
6)				No	0.07	0.793 to 0.863
7)				Yes	0.07	0.863 to 0.933

Random Number Generation

1st iteration \Rightarrow 0.43 (taking the 3rd record [which is incorrectly classified])

2nd iteration \Rightarrow 0.10 (taking the 2nd record)

3rd iteration \Rightarrow 0.73 (taking the 5th record)

: : :

Final Example

S.no	age	sex	chol	thalach	oldpeak	slope	TARGET
1	63	1	233	150	2.3	0	1
2	37	1	250	187	3.5	0	0
3	41	0	204	172	1.4	2	0
4	56	1	236	178	0.8	2	1
5	57	0	354	163	0.6	2	0
6	57	1	192	148	0.4	1	1
7	56	0	294	153	1.3	1	1

Initial Weight = $\frac{1}{n}$
to each observations

STEP 1

S.no	age	sex	chol	thalach	oldpeak	slope	TARGET	Weights
1	63	1	233	150	2.3	0	1	1/7
2	37	1	250	187	3.5	0	0	1/7
3	41	0	204	172	1.4	2	0	1/7
4	56	1	236	178	0.8	2	1	1/7
5	57	0	354	163	0.6	2	0	1/7
6	57	1	192	148	0.4	1	1	1/7
7	56	0	294	153	1.3	1	1	1/7

S.no	age	sex	chol	thalach	oldpeak	slope	TARGET	Weights
1	63	1	233	150	2.3	0	1	0.14
2	37	1	250	187	3.5	0	0	0.14
3	41	0	204	172	1.4	2	0	0.14
4	56	1	236	178	0.8	2	1	0.14
5	57	0	354	163	0.6	2	0	0.14
6	57	1	192	148	0.4	1	1	0.14
7	56	0	294	153	1.3	1	1	0.14

STEP 2

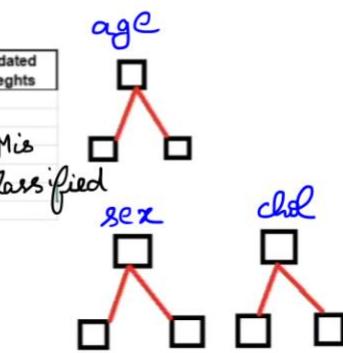
S.no	age	sex	chol	thalach	oldpeak	slope	TARGET	Weights	Predicted Output	Updated Weights
1	63	1	233	150	2.3	0	1	0.14	1	
2	37	1	250	187	3.5	0	0	0.14	0	
3	41	0	204	172	1.4	2	0	0.14	1	
4	56	1	236	178	0.8	2	1	0.14	1	
5	57	0	354	163	0.6	2	0	0.14	0	
6	57	1	192	148	0.4	1	1	0.14	0	
7	56	0	294	153	1.3	1	1	0.14	1	

$$\text{Total Error} = 2/7 = 0.28$$

Step 3 [based on age]

$$\text{Performance of stump} = \frac{1/2 \times \log_e(1 - \text{Total Error})}{\text{Total Error}}$$

$$\begin{aligned} \text{Stage-wise multiclass loss} &= \frac{1/2 \times \log_e(1 - 0.28)}{0.28} \\ &= 1/2 \times \log(2.5) \\ &= 1/2 \times .91 = 0.45 \end{aligned}$$



Age \Rightarrow having lowest Entropy

Correctly
classified

Incorrectly
classified

5 | 2

Old weights

S.no	age	sex	chol	thalach	oldpeak	slope	TARGET	Weights	Predicted Output	Updated Weights
1	63	1	233	150	2.3	0	1	0.14	1	
2	37	1	250	187	3.5	0	0	0.14	0	
3	41	0	204	172	1.4	2	0	0.14	1	
4	56	1	236	178	0.8	2	1	0.14	1	
5	57	0	354	163	0.6	2	0	0.14	0	
6	57	1	192	148	0.4	1	1	0.14	0	
7	56	0	294	153	1.3	1	1	0.14	1	

Performance of stump = 0.45

Step 4 : Updating Weights

New Weights = old Weights $\times e^{\pm \text{Performance}}$

+ for Misclassification
- for right classification

$$\text{New Weights} = 0.14 \times e^{+0.45} = 0.14 \times 1.56 = .21$$

$$\text{New Weights} = 0.14 \times e^{-0.45} = 0.14 \times 0.63 = 0.08$$

Probability of choosing correctly classified sample will be less

Probability of choosing this incorrectly classified sample will be more

For the upcoming decision trees

* Purposefully with the new weights formulae

Decreasing the probability
(correctly classified)
samples

To pay less attention for
the further decision tree

Increasing the probability
(incorrectly classified)
samples

To pay more attention for
the further decision tree

$$\text{New Weights} = 0.14 \times e^{+0.45} = 0.14 \times 1.56 = .21 \quad \text{For mis classification}$$

$$\text{New Weights} = 0.14 \times e^{-0.45} = 0.14 \times 0.63 = 0.08 \quad \text{For right classification}$$

✓

S.no	age	sex	chol	thalach	oldpeak	slope	TARGET	Weights	Predicted Output	Updated Wiegts
1	63	1	233	150	2.3	0	1	0.14	1	0.08
2	37	1	250	187	3.5	0	0	0.14	0	0.08
3	41	0	204	172	1.4	2	0	0.14	1	0.21
4	56	1	236	178	0.8	2	1	0.14	1	0.08
5	57	0	354	163	0.6	2	0	0.14	0	0.08
6	57	1	192	148	0.4	1	1	0.14	0	0.21
7	56	0	294	153	1.3	1	1	0.14	1	0.08
								1		0.82

Normal weights sum $\rightarrow 1$.

Updated weights sum is not equal to 1.

Step 5



Normalization

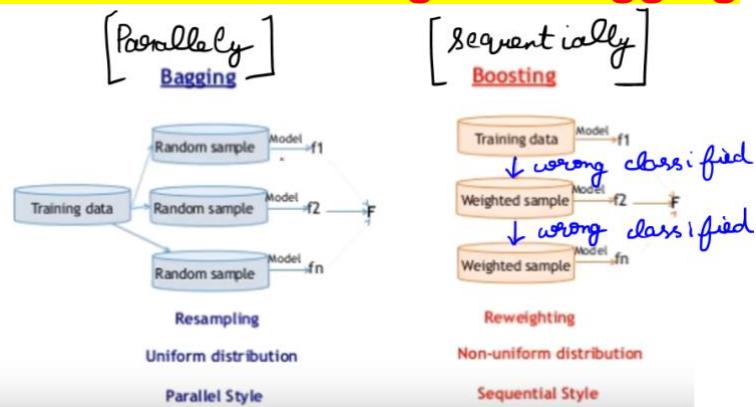
S.no	age	sex	chol	thalach	oldpeak	slope	TARGET	Weights	Predicted Output	Updated Wiegts	Normalized Weights
1	63	1	233	150	2.3	0	1	0.14	1	0.08	0.10
2	37	1	250	187	3.5	0	0	0.14	0	0.08	0.10
3	41	0	204	172	1.4	2	0	0.14	1	0.21	0.26
4	56	1	236	178	0.8	2	1	0.14	1	0.08	0.10
5	57	0	354	163	0.6	2	0	0.14	0	0.08	0.10
6	57	1	192	148	0.4	1	1	0.14	0	0.21	0.26
7	56	0	294	153	1.3	1	1	0.14	1	0.08	0.10
								1		0.82	1.00

Normalized weights

$$\text{Correctly classified sample} \Rightarrow \frac{0.08}{0.82} \Rightarrow 0.10$$

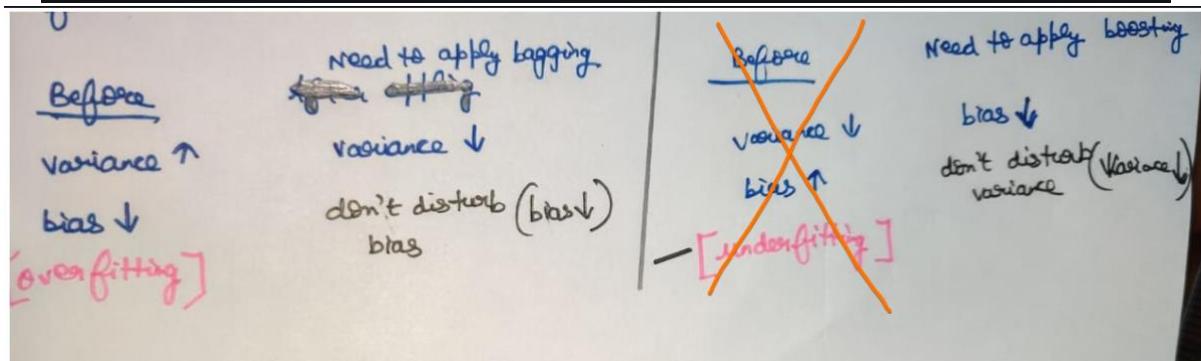
$$\text{In-correctly classified sample} \Rightarrow \frac{0.21}{0.82} \Rightarrow 0.26$$

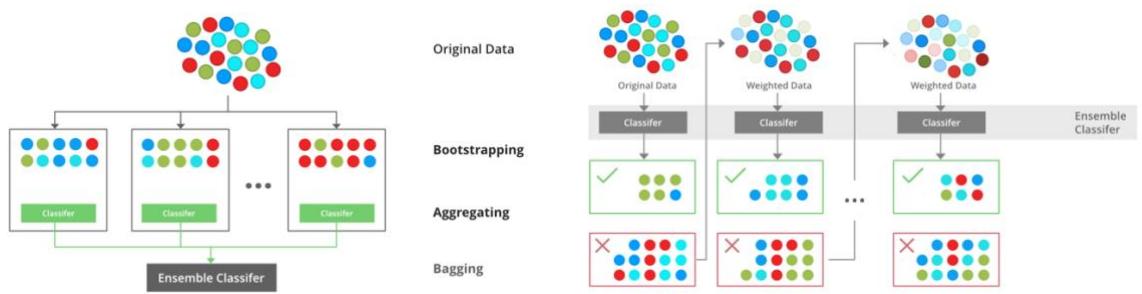
Difference between Boosting and Bagging



Differences Between Bagging and Boosting

S.NO	Bagging	Boosting
1.	The simplest way of combining predictions that belong to the same type.	A way of combining predictions that belong to the different types.
2.	Aim to decrease variance, not bias.	Aim to decrease bias, not variance.
3.	Each model receives equal weight.	Models are weighted according to their performance.
4.	Each model is built independently.	New models are influenced by the performance of previously built models.
5.	Different training data subsets are selected using row sampling with replacement and random sampling methods from the entire training dataset.	Every new subset contains the elements that were misclassified by previous models.
6.	Bagging tries to solve the over-fitting problem.	Boosting tries to reduce bias.
7.	If the classifier is unstable (high variance), then apply bagging.	If the classifier is stable and simple (high bias) the apply boosting.
8.	In this base classifiers are trained parallelly.	In this base classifiers are trained sequentially.
9	Example: The Random forest model uses Bagging.	Example: The AdaBoost uses Boosting techniques





Ada-Boost Sum

Season	Temperature	Play - Cricket	Predicted	Weight	New weight	Normalizing weights
1) Rain	Mild	Yes	Yes	$\frac{1}{8}$	0.047	0.042
2) Rain	Cool	Yes	Yes	$\frac{1}{8}$	0.047	0.012
3) Rain	Cool	No	Yes	$\frac{1}{8}$	0.047	0.012
4) Overcast	Cool	No	No	$\frac{1}{8}$	0.329	0.090
5) Sunny	Mild	Yes	Yes	$\frac{1}{8}$	0.047	0.012
6) Sunny	Cool	Yes	Yes	$\frac{1}{8}$	0.047	0.012
7) Rain	Mild	Yes	Yes	$\frac{1}{8}$	0.047	0.012
8) Rain	Mild	Yes	Yes	$\frac{1}{8}$	0.047	0.012

Step → Identify and assign equal weights to the model

Step → Identify a weak classifier in the model.

Step → Identify a weak classifier in the model
Assume 4 data has been weakly classified

Step → Calculate the total error

Set of 8 samples 4 sample ⇒ weakly classified 2/8

Total error ⇒ $\frac{1}{8}$

Step → Calculate the error in the model

$$E \Rightarrow \frac{1}{2} \log_e \left(\frac{1-E}{E} \right)$$

$$\Rightarrow \frac{1}{2} \log_e \left[\frac{1 - \frac{1}{8}}{\frac{1}{8}} \right]$$

$$E \Rightarrow 0.97$$

Step → Update the weights for strong classifier

$$\text{New weight} \Rightarrow \text{Old weight} * e^{-E}$$

$$\Rightarrow \frac{1}{8} * e^{-0.97}$$

$$\text{New weight} \Rightarrow 0.047$$

Step → Normalize the weights

Step → Update the weights for weak classifier

$$\text{New weight} \Rightarrow \text{Old weight} * e^E$$

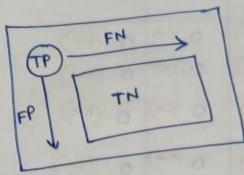
$$\Rightarrow \frac{1}{8} * e^{+0.97}$$

$$\text{New weight} \Rightarrow 0.329$$

$$\frac{0.047}{3.619} = 0.012 \quad \frac{0.329}{3.619} = 0.090$$

Calculate the Precision, Recall and Accuracy for below confusion matrix? What is the drawback of Accuracy metric?

	0	1	2	3	4	5	6	7	8	9
0	38	0	0	0	1	0	0	0	0	0
1	0	29	0	0	0	1	0	0	1	0
2	0	0	34	0	0	0	0	0	0	1
3	1	1	0	32	0	0	0	0	0	0
4	0	0	0	0	39	0	0	2	0	0
5	0	1	0	0	0	32	0	0	0	1
6	0	0	0	0	0	0	34	0	0	0
7	0	0	0	0	0	0	0	47	1	0
8	1	1	0	0	0	0	0	0	33	0
9	0	1	0	0	0	0	0	0	1	27



Feature - 0

$$TP \Rightarrow 38$$

$$FP \Rightarrow 2$$

$$TN \Rightarrow \cancel{(33+34+32+39+33+34+49+36+29)} \Rightarrow 319$$

$$FN \Rightarrow 1$$

$$TPR \Rightarrow \frac{TP}{TP+FN} \Rightarrow \frac{38}{38+1} \Rightarrow 0.97$$

$$FPR \Rightarrow \frac{FP}{FP+TN} \Rightarrow \frac{2}{2+319} \Rightarrow \frac{2}{321} \Rightarrow 0.00623$$

$$(0.00623, 0.97)$$

Feature - 5

$$TP \Rightarrow 32$$

$$FP \Rightarrow 0$$

$$TN \Rightarrow \cancel{(34+47+35+27)} \Rightarrow 143$$

$$FN \Rightarrow 1$$

$$TPR \Rightarrow \frac{TP}{TP+FN} \Rightarrow \frac{32}{32+1} \Rightarrow \frac{32}{33} \Rightarrow 0.96$$

$$FPR \Rightarrow \frac{FP}{FP+TN} \Rightarrow \frac{0}{0+143} \Rightarrow 0$$

$$(0, 0.96)$$

Stacking

References

<https://inblog.in/ENSEMBLE-METHODS-Bagging-Boosting-and-Stacking-WkkvEKBse>

<https://www.youtube.com/watch?v=Xz0x-8-cgaQ>

<https://www.analyticsvidhya.com/blog/2020/02/what-is-bootstrap-sampling-in-statistics-and-machine-learning/>

<https://www.youtube.com/watch?v=KIOeZ5cFZ50>

<https://www.youtube.com/watch?v=ERBr7jAcuvE>

<https://www.geeksforgeeks.org/bagging-vs-boosting-in-machine-learning/>