

Random Forest

We build Random Forest with multiple decision trees. If one classifier (Random forest here) is made of multiple classifiers (Decision tree), this method is called Ensemble Learning.

The concept is that, when multiple weak learners works together, it works better than one perfect learner. So, we will make multiple weak decision trees, whose combined decision will give us more accurate prediction.

Steps in Random Forest:

1. Create multiple bootstrapped data set.
 - Loop → for each bootstrapped dataset (Bagging)
(Create decision trees)
 - Loop → for each node in the decision tree
 - 2. Randomly select a few features
 - 3. compute Gini Impurity and decide the node
 - Store the tree
 - 4. Majority voting for classification decision.

Let us first build Random Forest on the given dataset.

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day 1	Sunny	Hot	High	Weak	NO NO
Day 2	Sunny	Hot	High	Strong	NO NO
Day 3	Ovencast	Hot	High	Weak	Yes
Day 4	Rain	Mild	High	Weak	Yes
Day 5	Rain	Cool	Normal	Weak	Yes
Day 6	Rain	Cool	Normal	Strong	NO
Day 7	Ovencast	Cool	Normal	Strong	Yes
Day 8	Sunny	Mild	High	Weak	NO
Day 9	Sunny	Cool	Normal	Weak	Yes
Day 10	Rain	Mild	Normal	Weak	Yes
Day 11	Sunny	Mild	Normal	Strong	Yes
Day 12	Ovencast	Mild	High	Strong	Yes
Day 13	Ovencast	Hot	Normal	Weak	Yes
Day 14	Rain	Mild	High	Strong	NO

① Bootstrapping

We will create multiple bootstrapped dataset at first. Bootstrapping is a method in statistics and ML whence a random number of data/^{on instances} gets selected. The word "Random" forest in random forest comes from this random selection of data/instances.

Let us create 3 bootstrapped dataset with 6 instances each →

Bootstrapped dataset 1 -

Let's assume Day 10, Day 11, Day 12, Day 13, Day 14 and Day 2 gets selected → copied →

Day	outlook	temperature	Humidity	Wind	Play Tennis
Day 10	Rain	Mild	Normal	Weak	Yes
Day 11	Sunny	Mild	Normal	Strong	Yes
Day 12	overcast	Mild	High	Strong	Yes
Day 13	overcast	Hot	Normal	Weak	Yes
Day 14	Rainy	Mild	High	Strong	No
Day 2	sunny	Hot	High	Strong	No

Bootstrapped dataset 2 -

We will randomly copy Day 2, Day 2, Day 3, Day 4, Day 5 and Day 2 Again so in Bootstrapped dataset, instances can be repeated.

Day	outlook	Temperature	Humidity	Wind	Play Tennis
Day 1	sunny	Hot	High	Weak	No
Day 2	Sunny	Hot	High	Strong	No
Day 3	overcast	Hot	High	Weak	Yes
Day 4	Rain	Mild	High	Weak	Yes
Day 5	Rain	Cool	Normal	Weak	Yes
Day 2	sunny	Hot	High	Strong	No

Bootstrapped Dataset 3 :-

we will select/copy Day 6, Day 7, Day 8, Day 9, Day 10, Day 13 randomly.

Day	outlook	Temperature	Humidity	Wind	Play Tennis
Day 6	Rain	Cool	Normal	Strong	NO
Day 7	overcast	Cool	Normal	Strong	Yes
Day 8	Sunny	Mild	High	weak	NO
Day 9	Sunny	Cool	Normal	Weak	Yes
Day 10	Rain	Mild	Normal	Weak	Yes
Day 13	overcast	Hot	Normal	Weak	Yes

Bagging :

Now we have to create Decision trees, for each Bootstrapped dataset. This technique is called Bagging, where we create several models on different random samples (bootstrapped) dataset.

We will now create full Decision Trees from all Bootstrapped dataset.
Let us create the first Decision Tree from the Bootstrapped dataset -1. For this at first we need to decide, what will be the root node.

② Feature Selection -

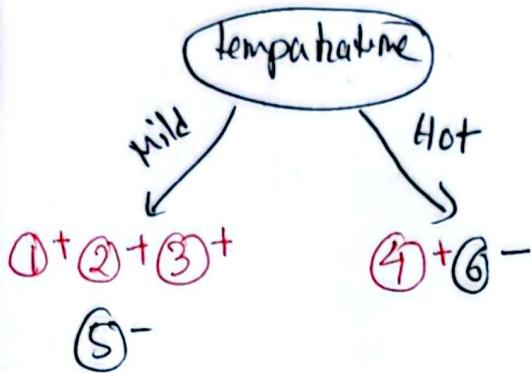
While creating the decision tree, whenever we are deciding a node, we will consider only a subset of features instead of all. Let us select Temperature and Humidity for root node from bootstrapped dataset 1.
So the training set will look like -

	Day	Temperature	Humidity	Play Tennis
①	Day 10	Mild	Normal	Yes
②	Day 11	Mild	Normal	Yes
③	Day 12	Mild	High	Yes
④	Day 13	Hot	Normal	Yes
⑤	Day 14	Mild	High	No
⑥	Day 2	Hot	High	No

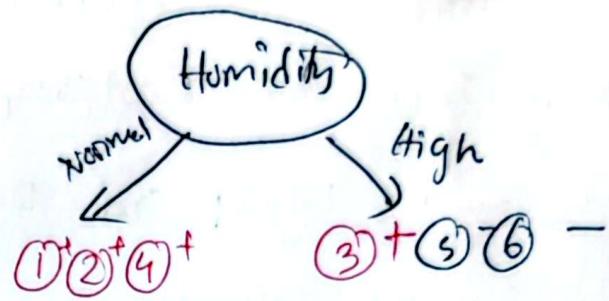
To create Decision Trees, we will use the Gini Impurity formula. Let us calculate Gini Impurity of Temperature and Humidity.

Calculating Gini Impurity:

Temperature



Humidity



$\Rightarrow \text{Gini}(\text{Temperature} = \text{Mild})$

$$= 1 - \left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 = 0.375$$

$\Rightarrow \text{Gini}(\text{Temperature} = \text{Hot})$

$$= 1 - \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 = 0.5$$

$\Rightarrow \text{Gini}(\text{Temperature})$

$$= \frac{4}{6} \times 0.375 + \frac{2}{6} \times 0.5$$

$$= 0.417$$

$\Rightarrow \text{Gini}(\text{Humidity} = \text{Normal})$

$$= 1 - \left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2 = 0$$

$\Rightarrow \text{Gini}(\text{Humidity} = \text{High})$

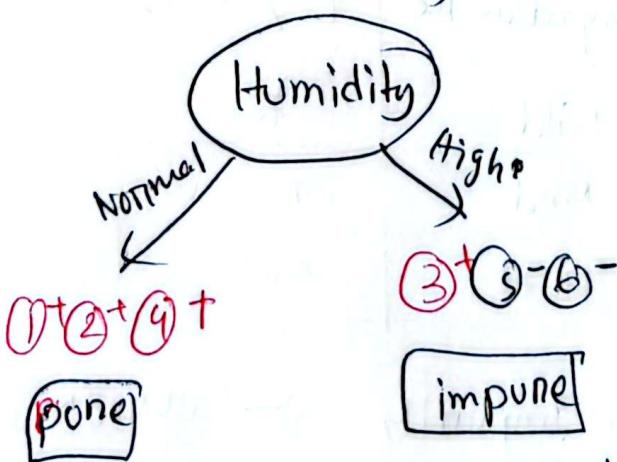
$$= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.44$$

$\Rightarrow \text{Gini}(\text{Humidity}) =$

$$= \frac{3}{6} \times 0 + \frac{3}{6} \times 0.44$$

$$= 0.222\overline{3}$$

As $\text{Gini}(\text{Temperature}) > \text{Gini}(\text{Humidity})$, we will consider Humidity as root node.



Decision = Yes

we need to again decide a feature for this node as split.

For the Impure Node, we have to do the process Again.
 - At first we will separate the datapoints of the impure node From the Bootstrapped Dataset-1.

~~Data~~ Day 12, Day 14 and Day 2 have ended up in the impure branch.

Data Set for the next level Node -

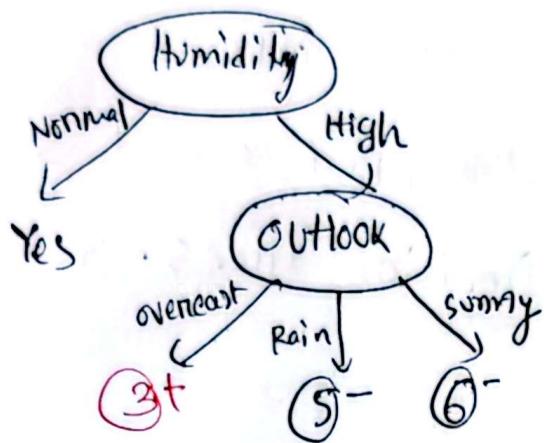
Day	outlook	temperature	Humidity	Wind	Play Tennis
Day 12	overcast	mild	high	strong	Yes
Day 14	Rain	mild	high	strong	No
Day 2	sunny	hot	high	strong	No

Again we will Select two features to decide the split for the Node. Let us select outlook and Temperature →

	Day	outlook	Temperature	Play Tennis
③	Day 12	overcast	mild	Yes
⑤	Day 14	Rain	mild	No
⑥	Day 2	sunny	hot	No.

Now we will calculate Gini Impurity of outlook and Temperature to decide the split.

OUTLOOK



$\text{Gini}(\text{outlook} = \text{overcast})$

$$= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.33$$

$\Rightarrow \text{Gini}(\text{outlook} = \text{Rain})$

$$= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.33$$

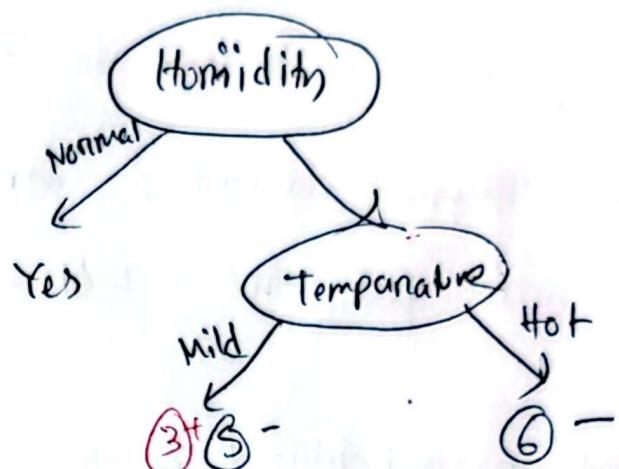
$\Rightarrow \text{Gini}(\text{outlook} = \text{Sunny})$

$$= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.33$$

$\therefore \text{Gini}(\text{overcast}) = \frac{1}{3} \times 0.33$

$$+ \frac{1}{3} \times 0.33 + \frac{1}{3} \times 0.33 \\ = 0$$

Temperature



$\Rightarrow \text{Gini}(\text{Temperature} = \text{Mild})$

$$= 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

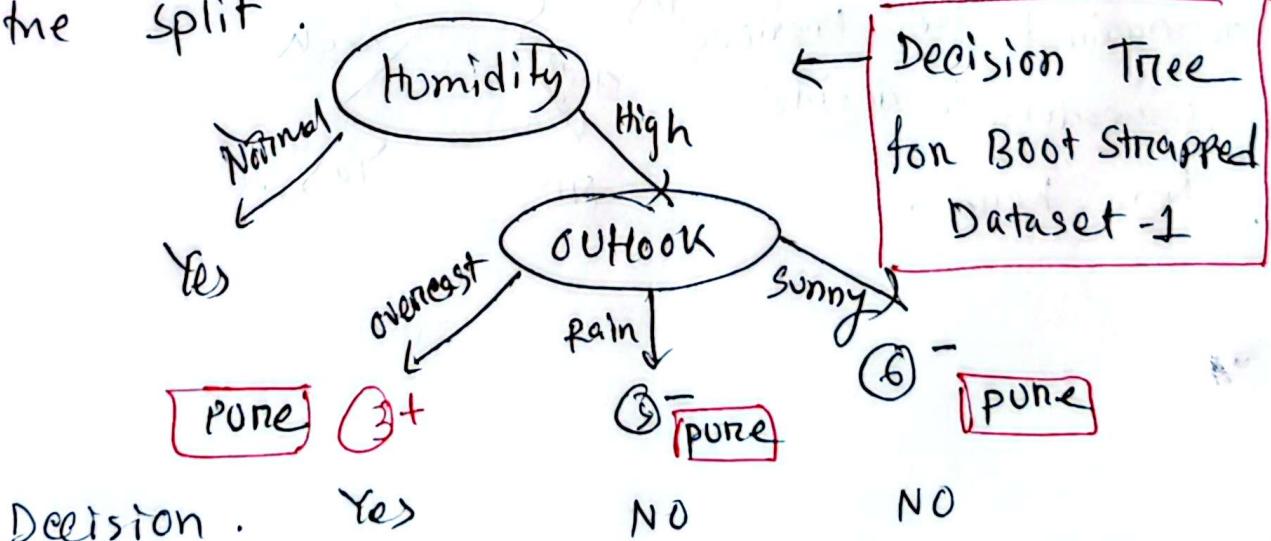
$\Rightarrow \text{Gini}(\text{Temperature} = \text{Hot})$

$$= 1 - \left(\frac{0}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0$$

$\therefore \text{Gini}(\text{Temperature}) = \frac{2}{3} \times 0.5$

$$+ \frac{1}{3} \times 0 = 0.33$$

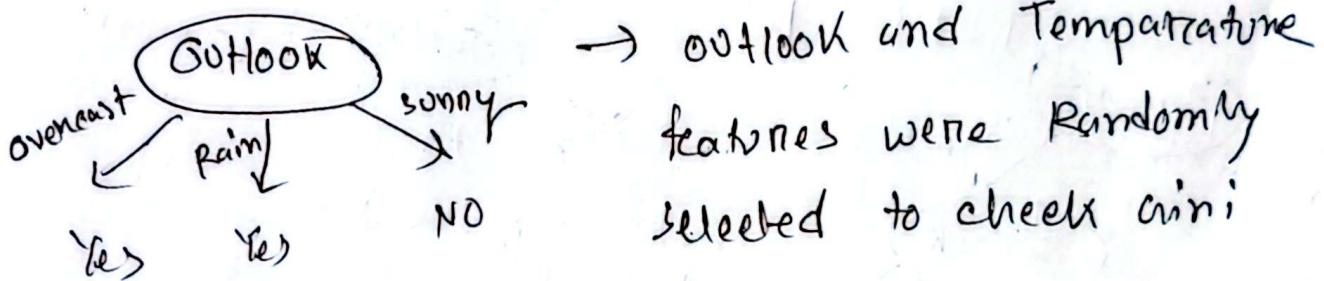
As $\text{Gini}(\text{outlook}) < \text{Gini}(\text{Temperature})$ we will use OUTLOOK as the split.



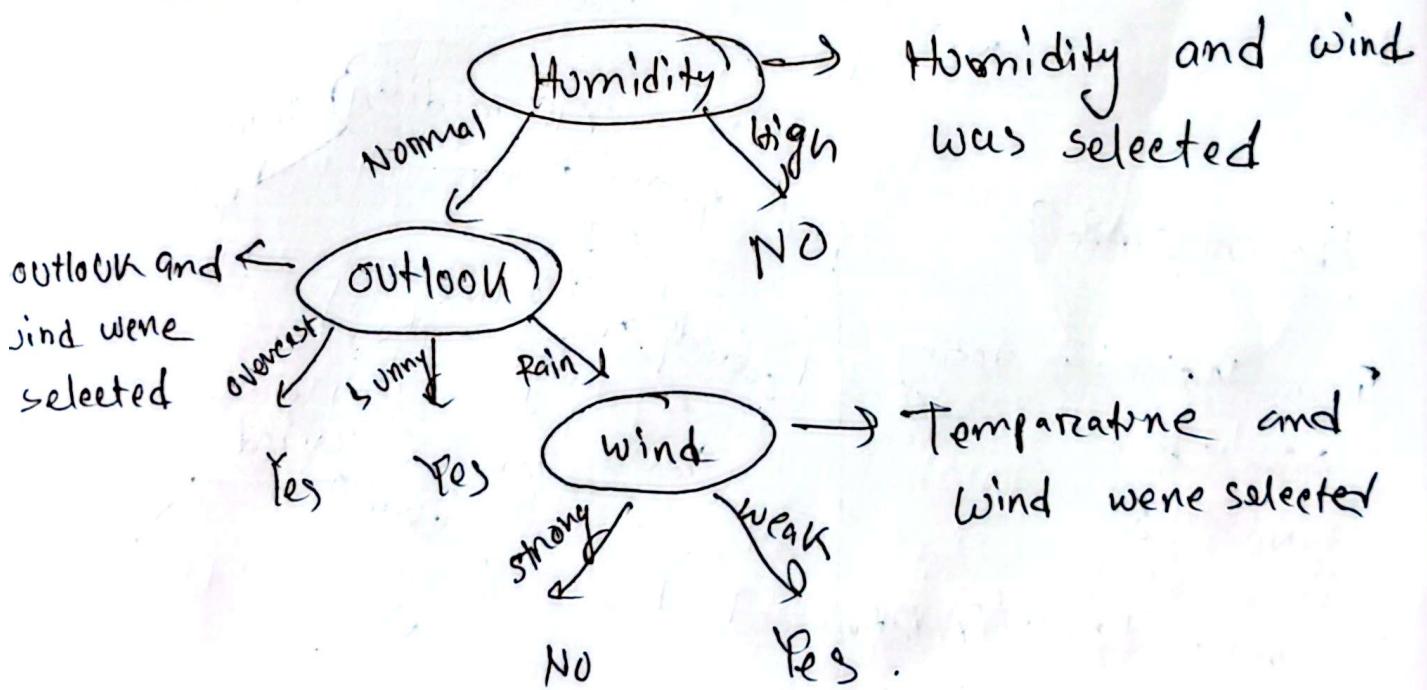
we will stop growing the tree as all leaves are pure.

Similarly if we find ~~dataset~~ ^{Decision Tree} for Bootstrapped dataset -1 and Bootstrapped dataset -2, we will get the following Decision Trees.

Bootstrapped Dataset - 2's Decision Tree



Bootstrapped Dataset - 3's Decision Tree -

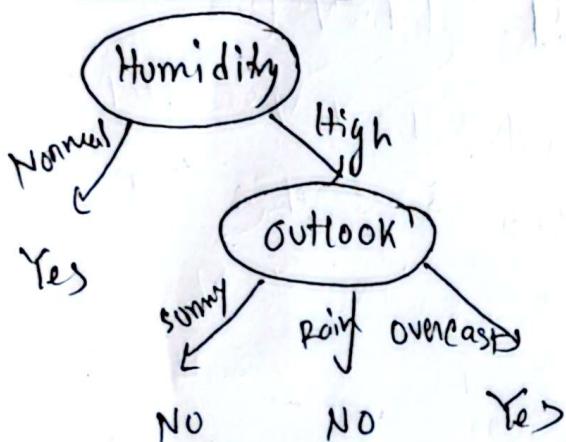


Now we have three Decision Trees of our Random Forest Model.

Let's predict ~~for~~ class for the following Query.

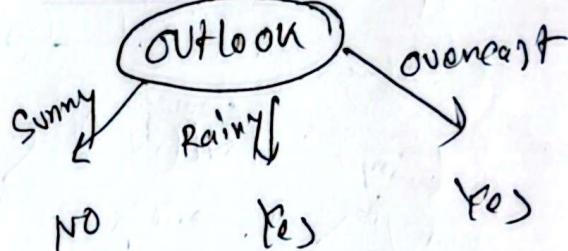
Day 13 →	<u>outlook</u>	<u>Temperature</u>	<u>Humidity</u>	<u>wind</u>	<u>Play Tennis</u>
	Rain	Hot	High	Weak	Yes

Decision Tree - 1 →



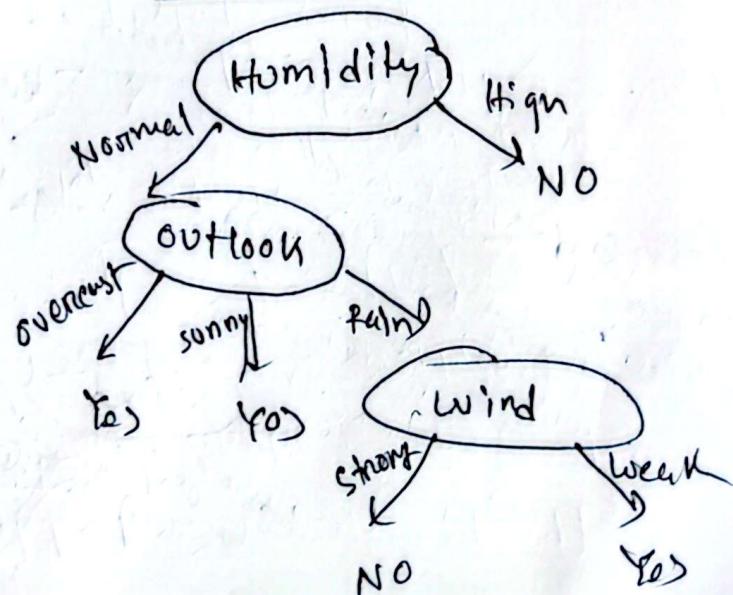
Query Result → NO.

Decision Tree - 2 →



Query Result → YES

Decision Tree - 3 →



Query Result → NO

From Majority Voting
YES: 2
NO: 2
The final Result/Prediction
of the Random Forest
Model is NO.