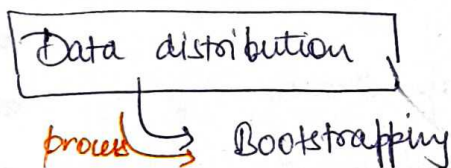
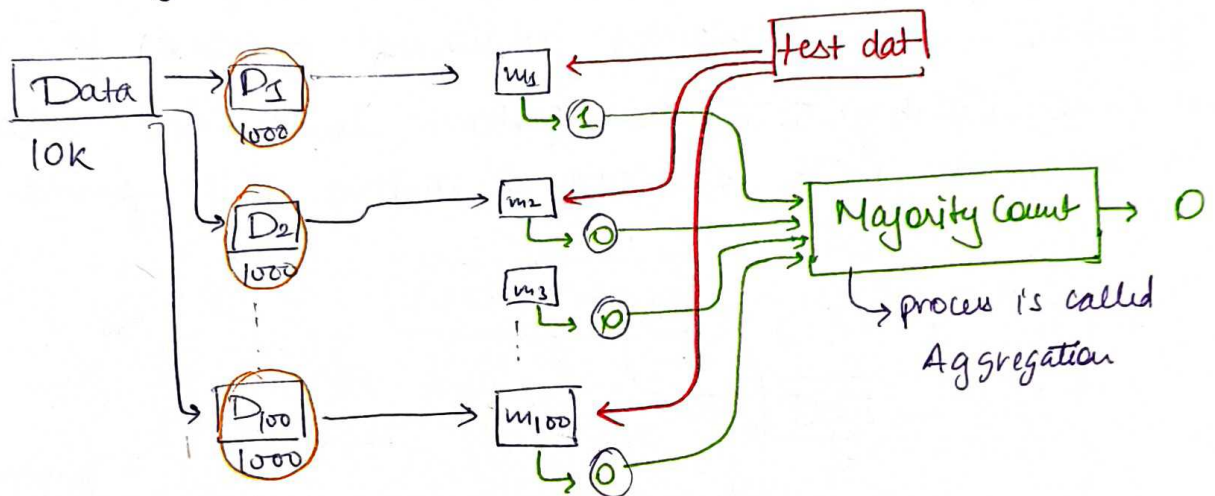
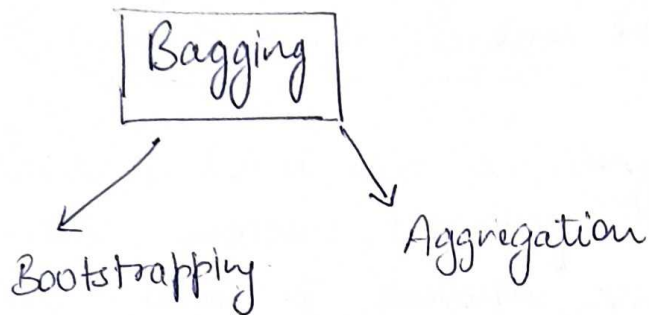
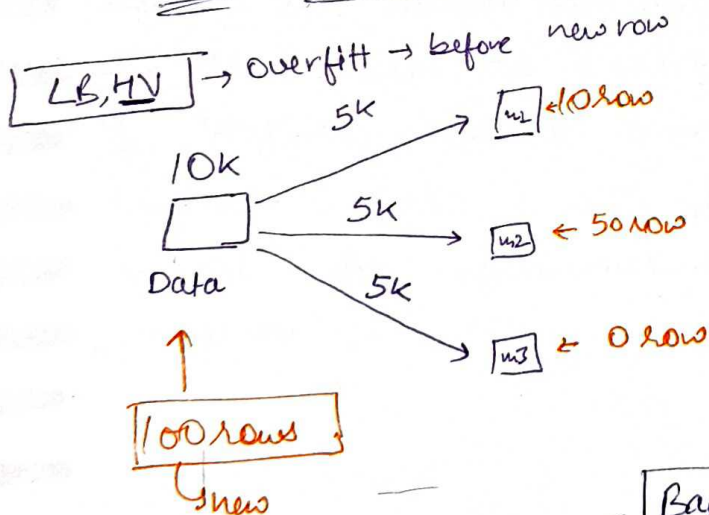


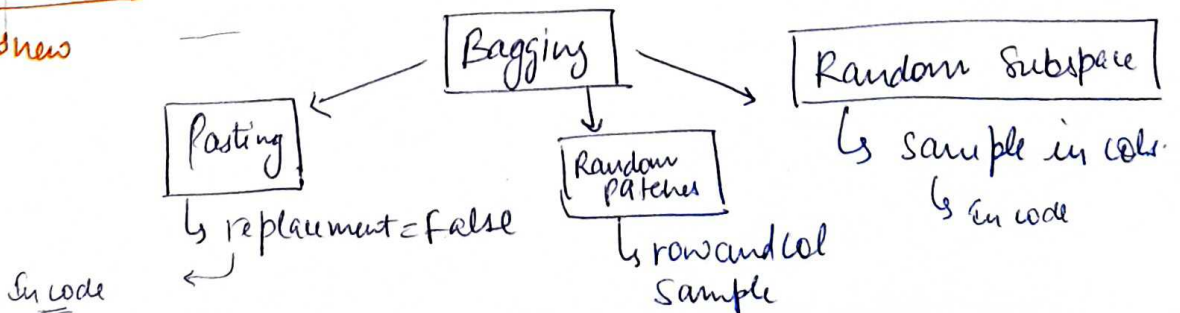
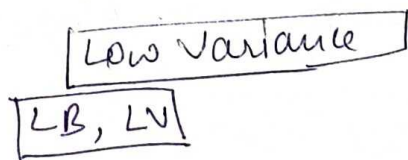
Bagging Ensemble



Bias Variance



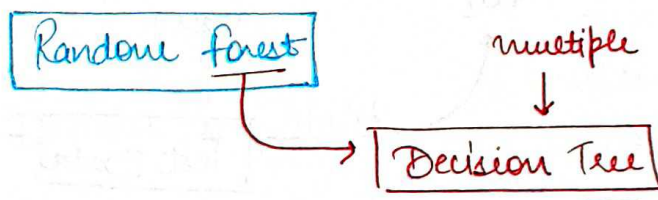
Not much changing in data. That's why doesn't affect on the model.



Random Forest

Introduction to Random Forest

Random Forest is a very versatile and widely used machine learning algorithm that belongs to the class of ensemble methods. Specifically, it is a type of bagging technique, which involves training many individual models (in this case, decision trees) and combining their outputs to make a final prediction.



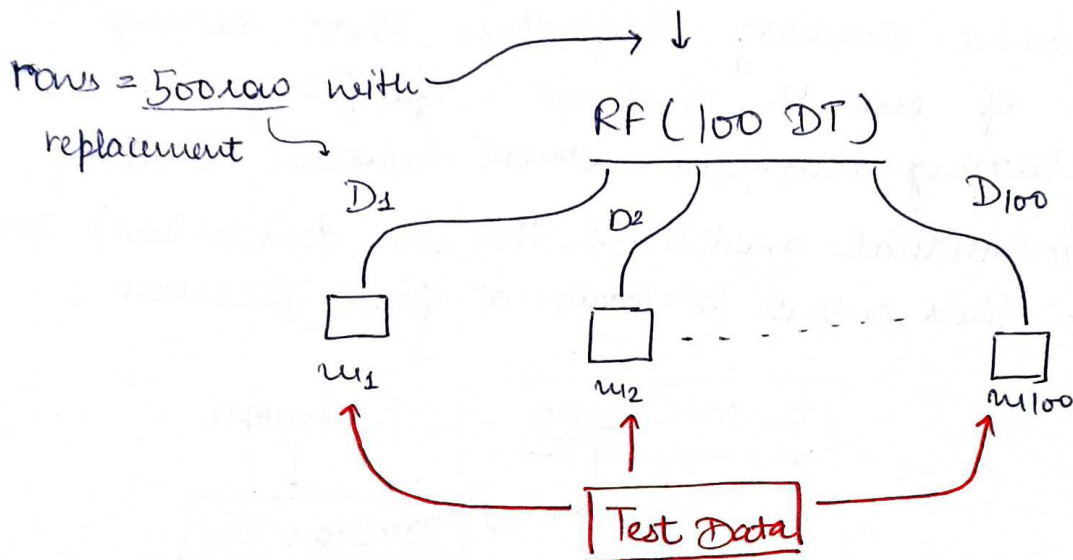
Bagging

Bagging, short for bootstrap aggregating, is a machine learning ensemble method designed to improve the stability and accuracy of machine learning algorithms used in statistical classification and regression. It also helps to avoid overfitting. The key principle of bagging is to generate multiple subsets of the original data (with replacement), train a separate model for each subset and then combine the results.

Random Forest Intuition

Random Forest \rightarrow base model \rightarrow Decision Tree

Data \rightarrow 100 rows / 5 cols (classification)



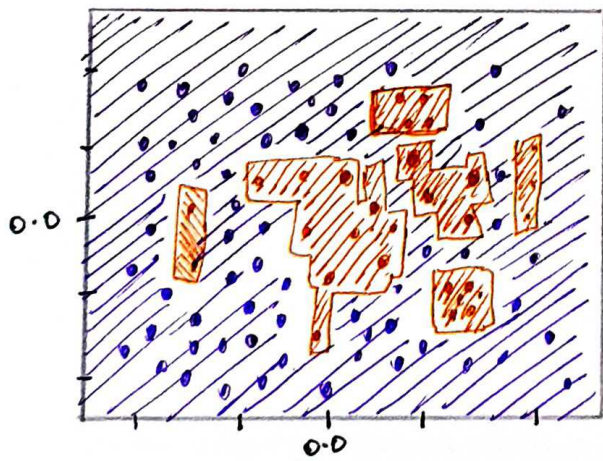
60 \rightarrow 0 } \rightarrow 0 majority count
40 \rightarrow 1

* In Regression, we find mean of the prediction accuracy of the model individually.

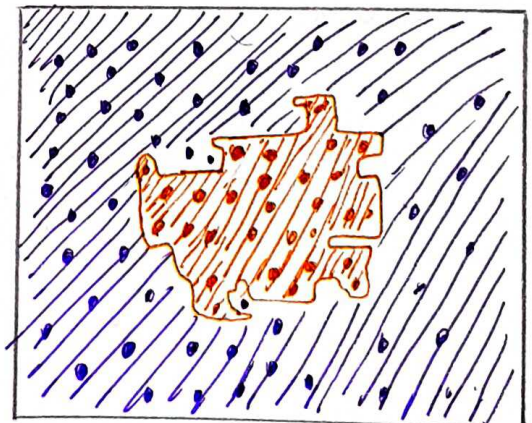
Why Random Forest Works well?

\downarrow

[Wisdom of the Crowd]
Decision Tree



Random Forest



Bias - Variable Tradeoff

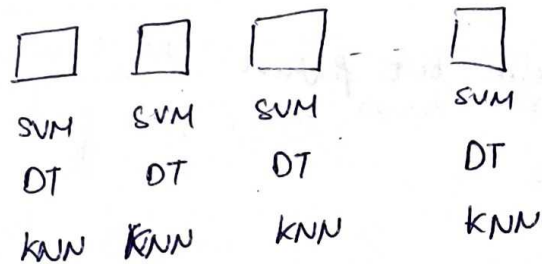
↓
Bias $\propto \frac{1}{\text{Variance}}$

LB, HV

↓
Random Forest → LB, LV

Bagging vs Random Forest

1) Bagging → any ML algorithm



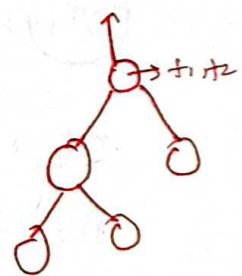
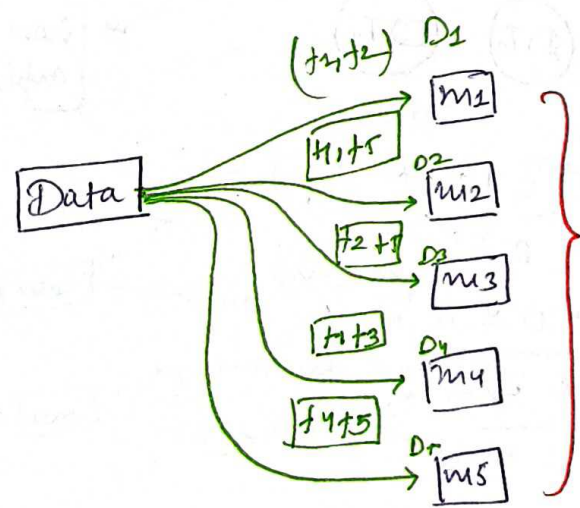
Random Forest → Decision Tree Model

Col samp = 50%

[feature decided before tree made]

2) Bagging →

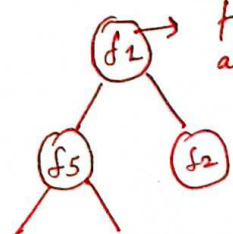
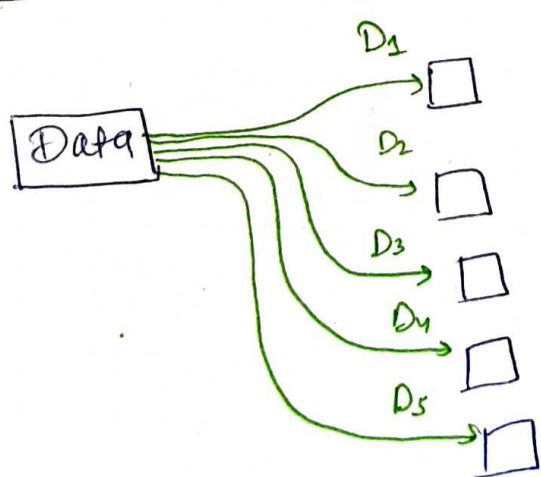
f1 | f2 | f3 | f4 | f5



Random Forest

Col sample = 50%

feature select at node level



* select cols 2 at node level.

Feature Importance

[Data] \rightarrow $\begin{bmatrix} f_1 & f_2 & f_3 & f_4 & f_5 \end{bmatrix}$

[RF]
 \downarrow

$f_1 \rightarrow 0.2$
 $f_2 \rightarrow 0.3$
 $f_3 \rightarrow 0.25$
 $f_4 \rightarrow 0.15$
 $f_5 \rightarrow 0.10$

} select best feature

(DT₁) (DT₂) (DT₃)

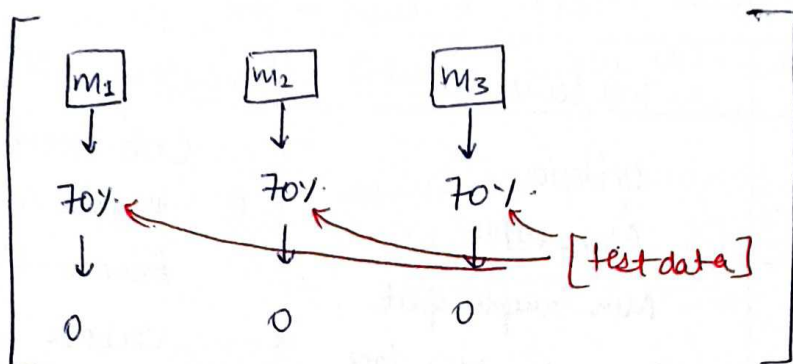
DT₁ $\rightarrow f_1 \rightarrow 0.1$

DT₂ $\rightarrow f_2 \rightarrow 0.3$

DT₃ $\rightarrow f_1 \rightarrow 0.4$
 $\underline{0.8/3}$

* how to find feature importance explained in decision tree

Why Ensemble Techniques work?



* So, we check whether the accuracy of ensemble learning is greater than single model ($EL > m_1$) ($EL > m_2$)

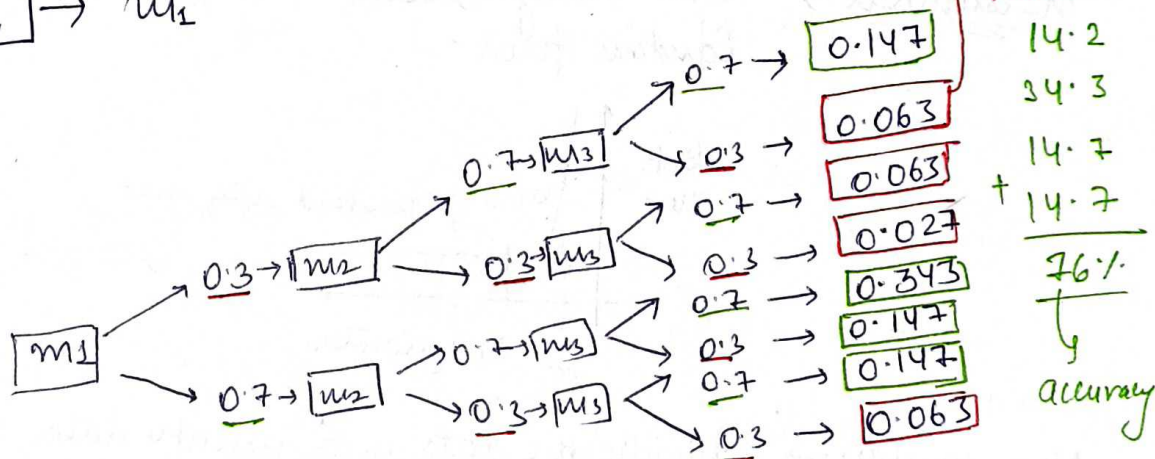
→ Assumption

m_1 , m_2 and m_3 are independent models

* If data is same but model is different

* If model is same but data is different.

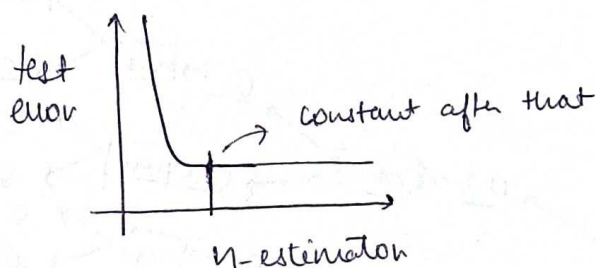
$x_q \rightarrow m_1$



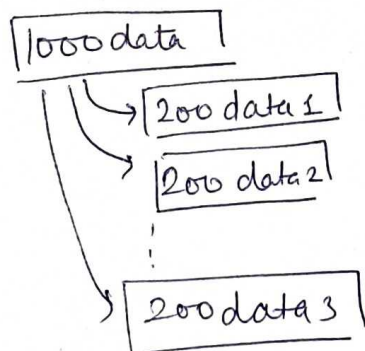
Random Forest Hyperparameters

| Forest level HP | Tree level HP | Miscellaneous HP |
|---|---|---|
| n-estimator Max features bootstrap Max-samples | criterion Max-depth Min-sample-split Min-samples-leaf Min-weight-fraction-leaf Max-leaf-nodes Min-impurity-decrease C _{yp} -alpha | Oob-score n-jobs Random-state Verbose Warm-start classweight |

n-estimator → how many Decision Tree train in side Random forest.



Max-samples → Divide the data into multiple data



How sample

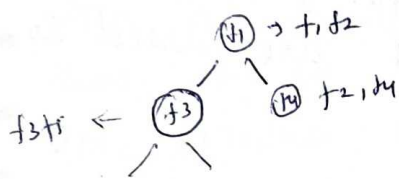
Bootstrap \rightarrow Row sampling with Replace or without replace

Bootstrap \rightarrow True \rightarrow repeat value in dataset

Bootstrap \rightarrow False \rightarrow Unique value in dataset

Max-features \rightarrow column sampling

5 columns \rightarrow



Max feature $\rightarrow \sqrt{5}$
No. of feature
 $\rightarrow \log(5)$

Warm-start \rightarrow training time is higher

True \rightarrow RF (n-est=10) \rightarrow 20 } not restart the process
if start the process where I stop
the process (11-20)

Class-weight \rightarrow imbalanced dataset

95% 5%
 \downarrow \downarrow
1 9

class weight give equal weight
both even data is imbalanced

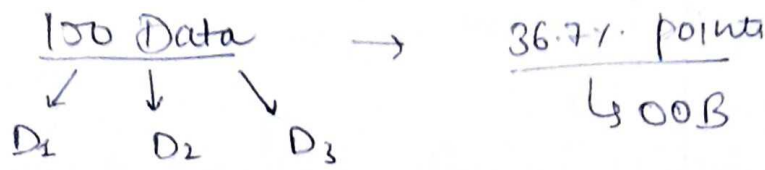
OOB Score

"OOB" stands for "out of bag". In the context of machine learning an out-of-bag score is a method of measuring the prediction error of random forests, bagging classifiers, and other ensemble methods that use bootstrap aggregation (bagging) where sub-samples of the training dataset are used to train individual models.

Here's how it works:

- 1.) Each tree in the ensemble is trained on a distinct bootstrap sampling, some samples from the dataset will be left out during the training of each tree. These samples are called "out-of-bag" samples.
- 2.) The out of bag samples can then be used as a validation set. We can pass them through the tree that didn't see them during training and obtain predictions.
- 3.) These predictions are then compared to the actual values to compute an "out-of-bag" score, which can be thought of as an estimate of the prediction error on unseen data.

One of the advantages of the out-of-bag score is that it allows us to estimate the prediction error without needing a separate validation set. This can be particularly useful when the dataset is small and partitioning it into training and validation set might leave too few samples for effective learning.



Extremely Randomized Trees

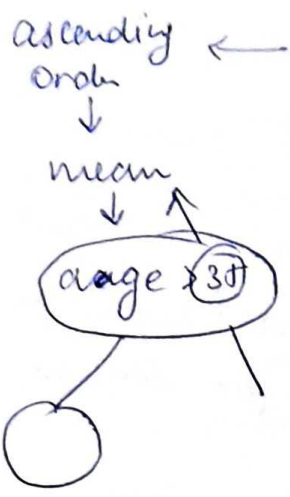
Extra Trees is short for "Extremely Randomized Trees". It's a modification of the Random Forest algorithm that changes the way the splitting points for decision tree branches are chosen.

In traditional decision tree algorithms (and therefore in Random Forests), the optimal split point for each feature is calculated, which involves a degree of computation. For a given node, the feature and the corresponding optimal split point that provide the best split are chosen. On the other hand, in the Extra Trees algorithm, for each feature under consideration, a split point is chosen completely at random. The best performing feature and its associated randomness to the model, hence the name "Extremely Randomized Trees".

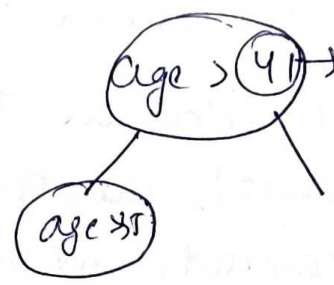
Because of this difference, Extra Trees tend to have more branches (be deeper) than Random Forests, and the splits are made arbitrarily. This can sometimes lead to models that perform better, especially on tasks where the data may not have clear optimal split points. However, like all models, whether Extra Trees will outperform Random Forests (or any other algo) depends on the specific dataset and task.

| age | insurance |
|-----|-----------|
| 27 | N |
| 43 | Y |

Decision Tree



Extremely Randomised Trees



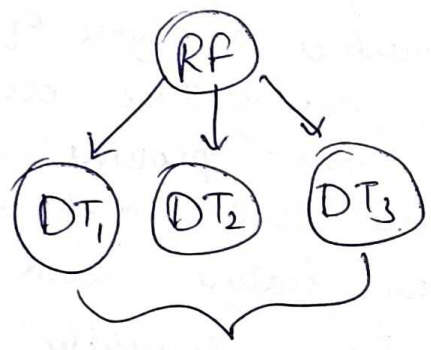
Random

high bias

training accuracy ↓

accuracy change

Extra Trees
Extra Randomness



independent
(decorrelation)

* Sometimes good work

Advantages

Robustness to Overfitting: Random Forests are less prone to overfitting compared to individual decision trees, because they average the results from many different trees, each of which might overfit the data in a different way.

Handling Large Datasets: They can handle large datasets with high dimensionality effectively.

Less pre-processing: Random Forest can handle both categorical and numerical variables without the need for scaling or normalization. They can also handle missing values.

Variable Importance: They provide insights into which features are most important in the prediction.

Parallelizable: The training of individual tree can be parallelized, as they are independent of each other. This speeds up the training process.

Non-Parametric: Random Forests are non-parametric, meaning they make no assumptions about the function form of the transformation from inputs to output. This makes them very flexible and able to model complex, non-linear relationships.

Disadvantages

- Model Interpretability: One of the biggest drawbacks of Random Forest is that they lack the interpretability of simpler models like linear regression or decision trees. While you can rank features by their importance, the model as a whole is essentially a black box.
- Performance with Unbalanced Data: Random Forest can be biased towards the majority class when dealing with unbalanced datasets. This can sometimes be mitigated by balancing the dataset prior to training.
- Predictive Performance: Although Random Forests generally perform well, they may not always provide the best predictive performance. Gradient Boosting machine for instance, often outperform Random Forests. If the ~~data~~ relationship within the data are linear, a linear model will likely perform better than a Random Forests.
- Inefficiency with Sparse Data: Random Forests might not be the best choice for sparse data or text data where linear model or other algo might be more suitable.

- Parameter Tuning: Although Random Forests requires less tuning than some other models, there are still several parameters (like the number of trees, tree depth etc) that can affect model performance and need to be optimized.
- Difficulty with High Cardinality Features: Random Forests can struggle with high cardinality categorical features (features with a large number of distinct values). These types of features can lead to trees that are biased towards the variable with more levels and may cause overfitting.
- Can't Extrapolate:- This is because they do not predict beyond the range of the training data, and that they may not predict as accurately as other regression models.

Extrapolation is not possible