

PLTHIYAGU ([HTTPS://ID.ANALYTICSVIDHYA.COM/ACCOUNTS/PROFILE/](https://id.analyticsvidhya.com/accounts/profile/))**Analytics Vidhya**

Learn everything about analytics

([https://www.analyticsvidhya.com/myfeed/?utm-](https://www.analyticsvidhya.com/myfeed/?utm-source=blog&utm-medium=top-icon/)[source=blog&utm-medium=top-icon/](https://www.analyticsvidhya.com/myfeed/?utm-source=blog&utm-medium=top-icon/))**Certified Program
AI & ML BlackBelt+****Reserve Your Seat**

- + 12+ Courses in AI, ML & DL
- + 25+ Latest Projects
- + Interview Guidance
- + 1:1 Mentorship Sessions

([https://courses.analyticsvidhya.com/bundles/certified-ai-ml-blackbelt-plus/?](https://courses.analyticsvidhya.com/bundles/certified-ai-ml-blackbelt-plus/?utm_source=blog&utm_medium=flash_strip&utm_campaign=bbplus_reserve_seat)[utm_source=blog&utm_medium=flash_strip&utm_campaign=bbplus_reserve_seat](https://courses.analyticsvidhya.com/bundles/certified-ai-ml-blackbelt-plus/?utm_source=blog&utm_medium=flash_strip&utm_campaign=bbplus_reserve_seat))[ALGORITHM \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/ALGORITHM/\)](https://www.analyticsvidhya.com/blog/category/algorithm/)[BEGINNER \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/BEGINNER/\)](https://www.analyticsvidhya.com/blog/category/beginner/)[BIAS AND VARIANCE \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/BIAS-AND-VARIANCE/\)](https://www.analyticsvidhya.com/blog/category/bias-and-variance/)[CLASSIFICATION \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/CLASSIFICATION/\)](https://www.analyticsvidhya.com/blog/category/classification/)[DATA SCIENCE \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/DATA-SCIENCE/\)](https://www.analyticsvidhya.com/blog/category/data-science/)[DATA VISUALIZATION \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/DATA-VISUALIZATION/\)](https://www.analyticsvidhya.com/blog/category/data-visualization/)

A Beginner's Guide to Random Forest Hyperparameter Tuning

SHAROON SAXENA ([HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/AUTHOR/SHAROON-SAXENA/](https://www.analyticsvidhya.com/blog/author/sharoon-saxena/)), MARCH 12, 2020

Industrial Building Rated AA - 156,627 SF Industrial I

Industrial Building For Lease | Food Quality Ready | Excellent Layout | 40' x 40' Column
30corporate.com

Introduction to Random Forest

What's the first image that comes to your mind when you think about Random Forest? It conjures up images of trees and a mystical and magical land. And that's what the Random Forest algorithm does!

It is an ensemble algorithm (https://www.analyticsvidhya.com/blog/2018/06/comprehensive-guide-for-ensemble-models/?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning), that combines multiple decision trees and navigates complex problems to give us the final result.



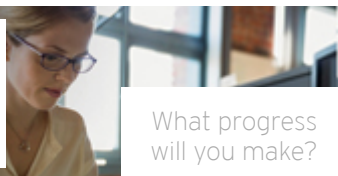
(https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/articles_img_main.jpg).

I've lost count of the number of times I've relied on the Random Forest algorithm in my machine learning projects and even hackathons. What makes random forest different from other ensemble algorithms is the fact that each individual tree is built on a subset of data and features.

Random Forest comes with a caveat – the numerous hyperparameters that can make fresher data scientists weak in the knees. But don't worry! In this article, we will be looking at the various Random Forest hyperparameters and understand how to tune and optimize them.

I assume you have a basic understanding of the random forest algorithm (and decision trees). If not, I encourage you to go through the below resources first:

*Compliance careers can
have a global impact.*



What progress
will you make?

COMPL
OPPOF
jobs.cil



- [Introduction to Random Forest \(https://www.analyticsvidhya.com/blog/2014/06/introduction-random-forest-simplified/?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning\)](https://www.analyticsvidhya.com/blog/2014/06/introduction-random-forest-simplified/?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning).
- [Getting Started with Decision Trees \(https://courses.analyticsvidhya.com/courses/getting-started-with-decision-trees?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning\)](https://courses.analyticsvidhya.com/courses/getting-started-with-decision-trees?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning).
(Free Course)

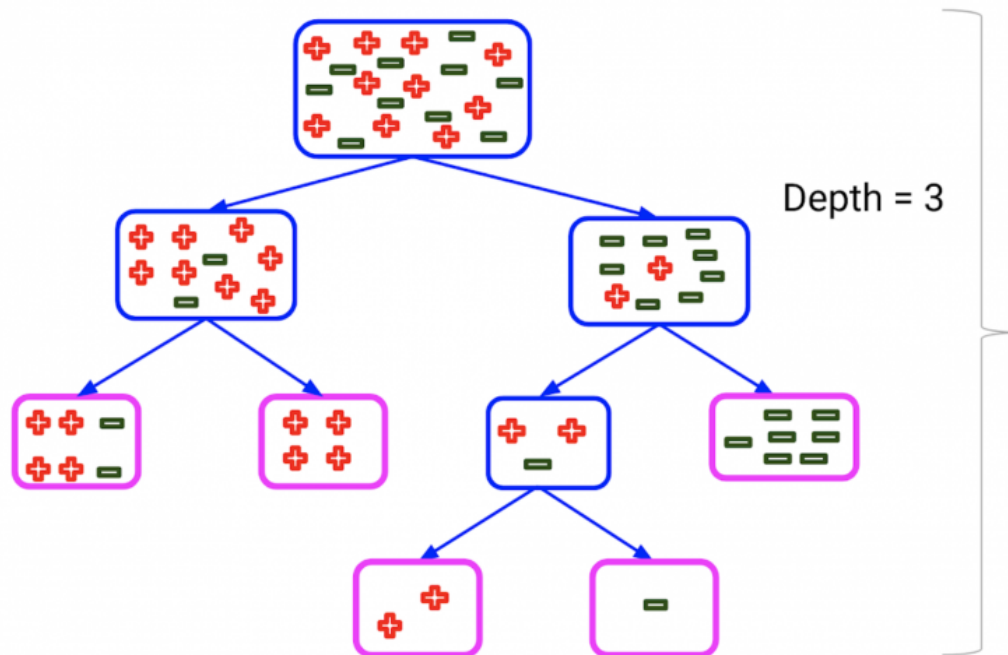
Random Forest Hyperparameters we'll be Looking at:

- max_depth
- min_sample_split
- max_leaf_nodes
- min_samples_leaf
- n_estimators
- max_sample (bootstrap sample)
- max_features

Random Forest Hyperparameter #1: max_depth

Let's discuss the critical *max_depth* hyperparameter first. The *max_depth* of a tree in Random Forest is defined as the longest path between the root node and the leaf node:

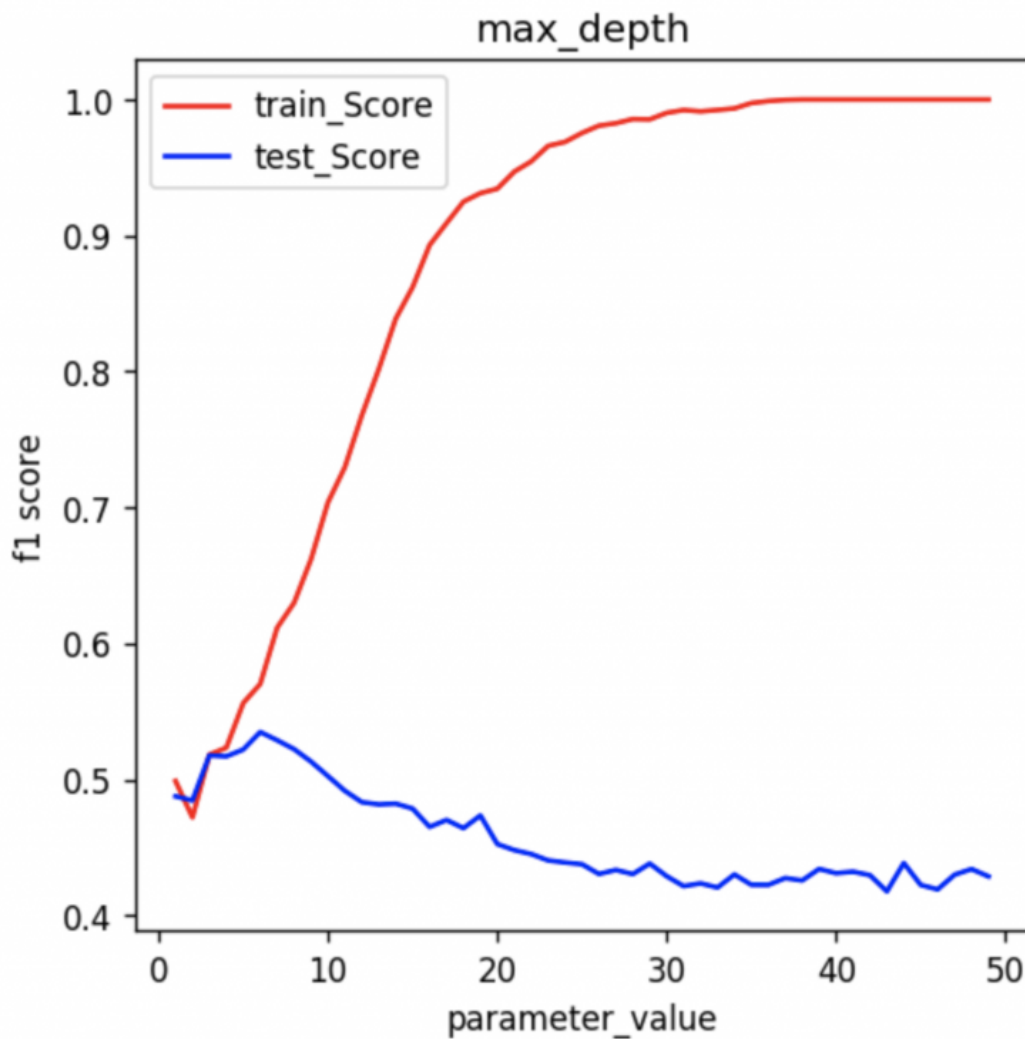





<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-15.07.03.png>

Using the *max_depth* parameter, I can limit up to what depth I want every tree in my random forest to grow.





(<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-15.08.50.png>).



All-In-One HR Soft

Intuitive HR Apps. Simplify HR Management, Reduce Busywork, & Money. Free Trial.

Learn More



In this graph, we can clearly see that as the max depth of the decision tree increases, the performance of the model over the training set increases continuously. On the other hand as the *max_depth* value increases, the performance over the test set increases initially but after a certain point, it starts to decrease rapidly.

Can you think of a reason for this? The tree starts to overfit the training set and therefore is not able to generalize over the unseen points in the test set.

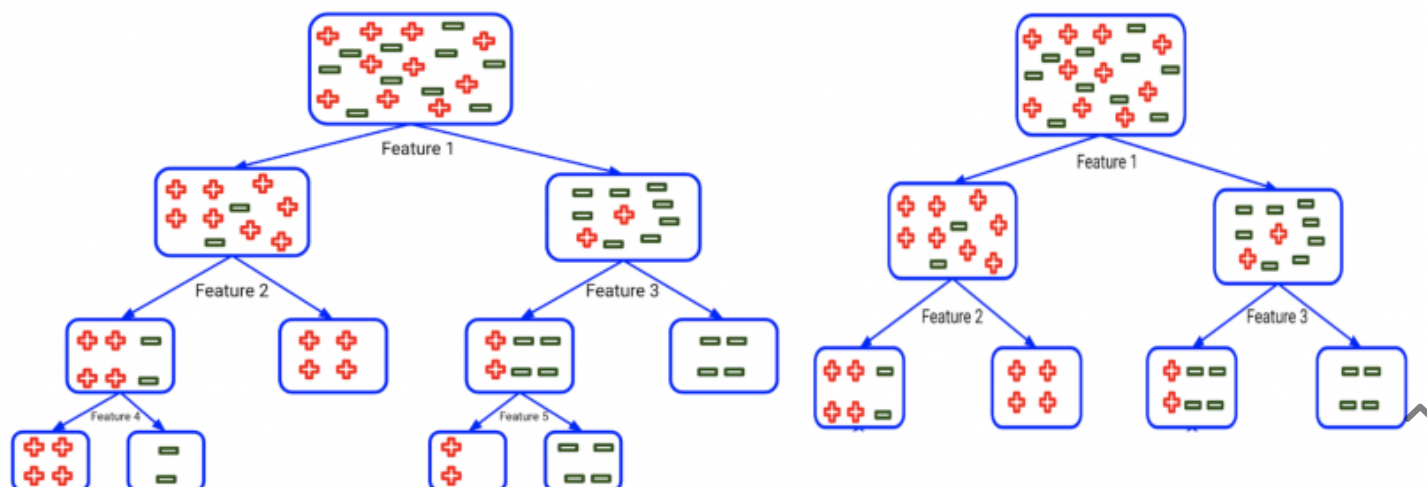
Among the parameters of a decision tree (https://www.analyticsvidhya.com/blog/2016/04/tree-based-algorithms-complete-tutorial-scratch-in-python/?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning), *max_depth* works on the macro level by greatly reducing the growth of the Decision Tree.

Random Forest Hyperparameter #2: min_sample_split

min_sample_split – a parameter that tells the decision tree in a random forest the minimum required number of observations in any given node in order to split it.

The default value of the minimum_sample_split is assigned to 2. This means that if any terminal node has more than two observations and is not a pure node, we can split it further into subnodes.

Having a default value as 2 poses the issue that a tree often keeps on splitting until the nodes are completely pure. As a result, the tree grows in size and therefore overfits the data.



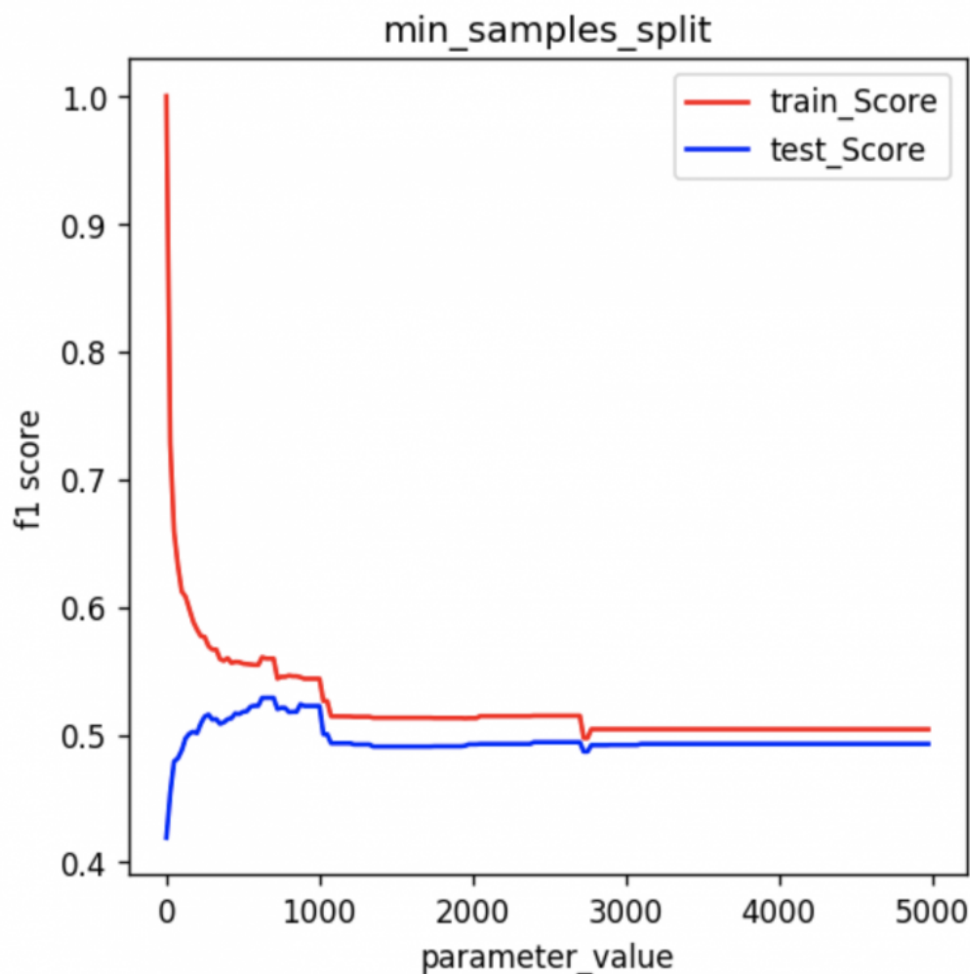
(<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-15.11.56.png>).

Industrial Building Rated AA - 156,627 SF Industrial I

Industrial Building For Lease | Food Quality Ready | Excellent Layout | 40' x 40' Column
30corporate.com

By increasing the value of the *min_sample_split*, we can reduce the number of splits that happen in the decision tree and therefore prevent the model from overfitting. In the above example, if we increase the *min_sample_split* value from 2 to 6, the tree on the left would then look like the tree on the right.

Now, let's look at the effect of *min_samples_split* on the performance of the model. The graph below is plotted considering that all the other parameters remain the same and only the value of *min_samples_split* is changed:



(<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-15.13.06.png>).

On increasing the value of the *min_sample_split* hyperparameter, we can clearly see that for the small value of parameters, there is a significant difference between the training score and the test scores. But as the value of the parameter increases, the difference between the train score and the test score decreases.

But there's one thing you should keep in mind. *When the parameter value increases too much, there is an overall dip in both the training score and test scores. This is due to the fact that the minimum requirement of splitting a node is so high that there are no significant splits observed. As a result, the random forest starts to underfit.*

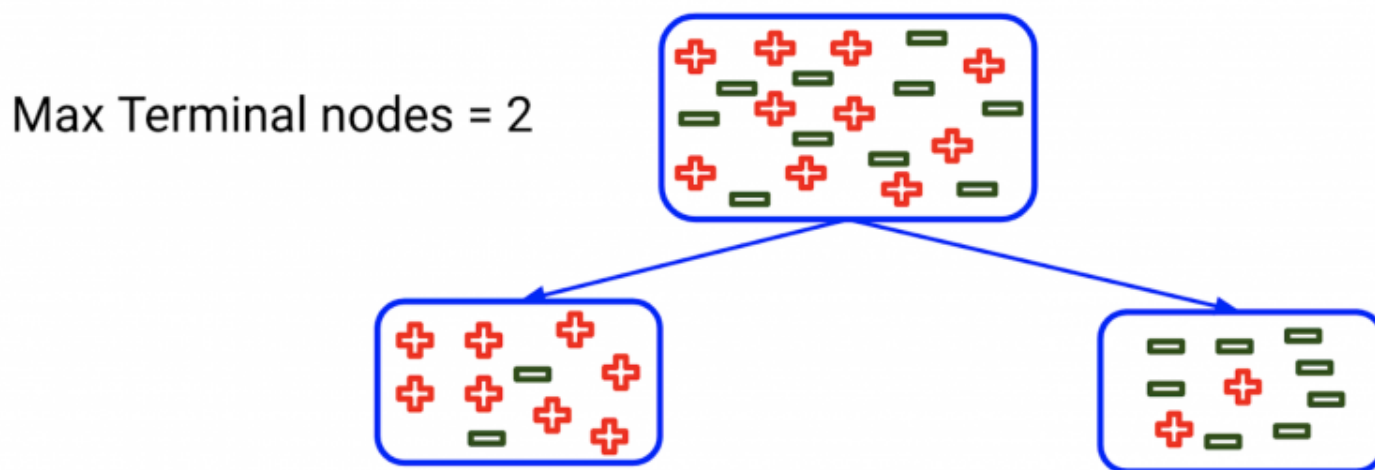
You can read more about the concept of overfitting and underfitting here:

- [Underfitting vs. Overfitting in Machine Learning](https://www.analyticsvidhya.com/blog/2020/02/underfitting-overfitting-best-fitting-machine-learning/?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning)
(https://www.analyticsvidhya.com/blog/2020/02/underfitting-overfitting-best-fitting-machine-learning/?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning)

Random Forest Hyperparameter #3: *max_terminal_nodes*

Next, let's move on to another Random Forest hyperparameter called *max_leaf_nodes*. **This hyperparameter sets a condition on the splitting of the nodes in the tree and hence restricts the growth of the tree.** If after splitting we have more terminal nodes than the specified number of terminal nodes, it will stop the splitting and the tree will not grow further.

Let's say we set the maximum terminal nodes as 2 in this case. As there is only one node, it will allow the tree to grow further:

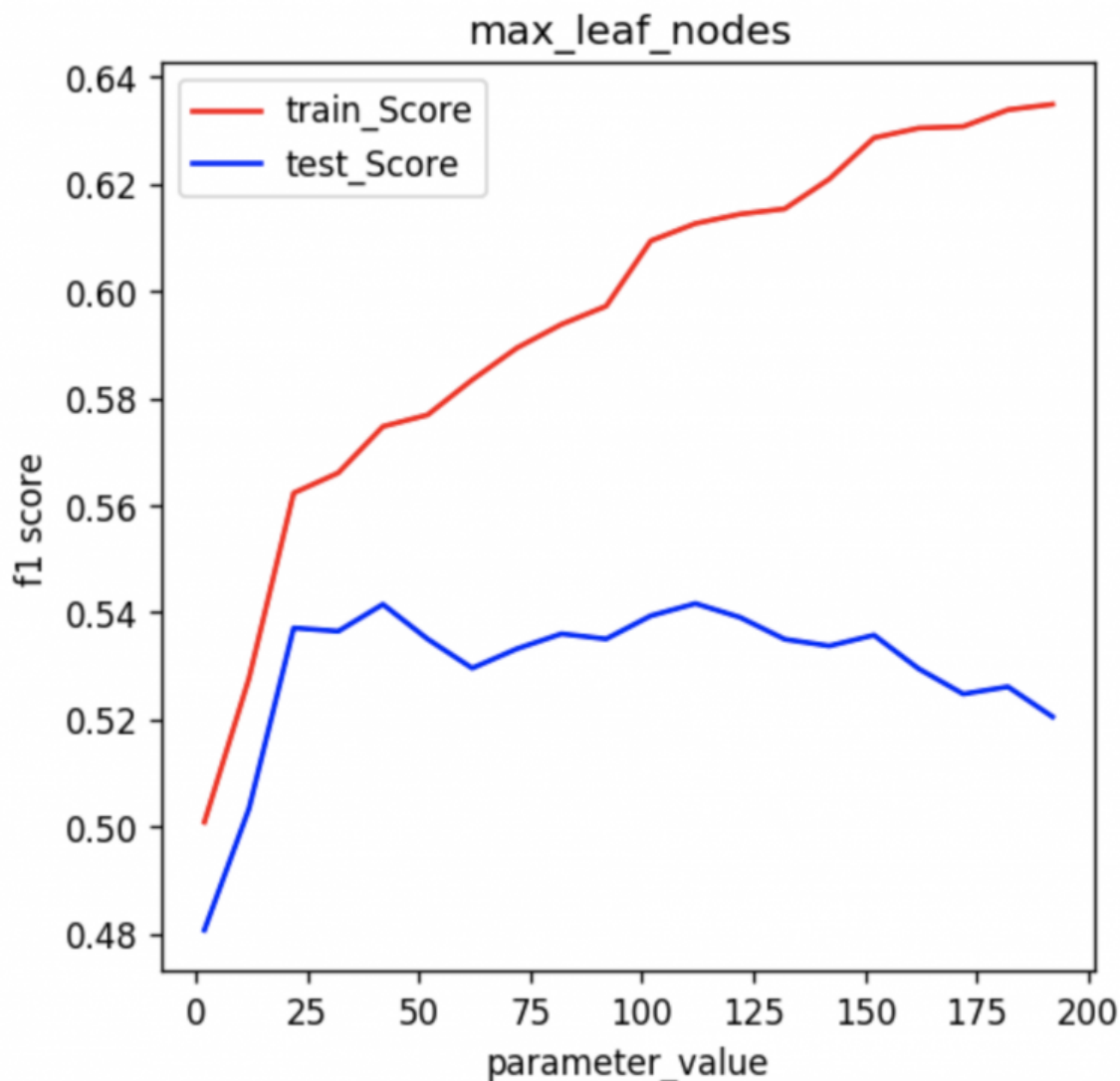


(<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-15.18.39.png>).



Now, after the first split, you can see that there are 2 nodes here and we have set the maximum terminal nodes as 2. Hence, the tree will terminate here and will not grow further. This is how setting the maximum terminal nodes or *max_leaf_nodes* can help us in preventing overfitting.

Note that if the value of the *max_leaf_nodes* is very small, the random forest is likely to underfit. Let's see how this parameter affects the random forest model's performance:



(<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-15.21.07.png>).

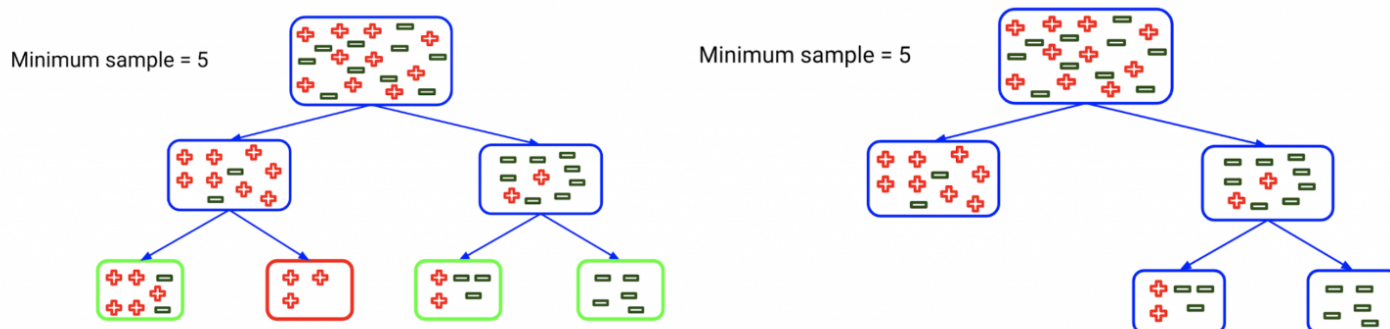
We can see that when the parameter value is very small, the tree is underfitting and as the parameter value increases, the performance of the tree over both test and train increases. According to this plot, the tree starts to overfit as the parameter value goes beyond 25.



Random Forest Hyperparameter #4: `min_samples_leaf`

Time to shift our focus to *min_sample_leaf*. This Random Forest hyperparameter specifies the minimum number of samples that should be present in the leaf node **after splitting** a node.

Let's understand *min_sample_leaf* using an example. Let's say we have set the minimum samples for a terminal node as 5:



(<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-15.34.46.png>).

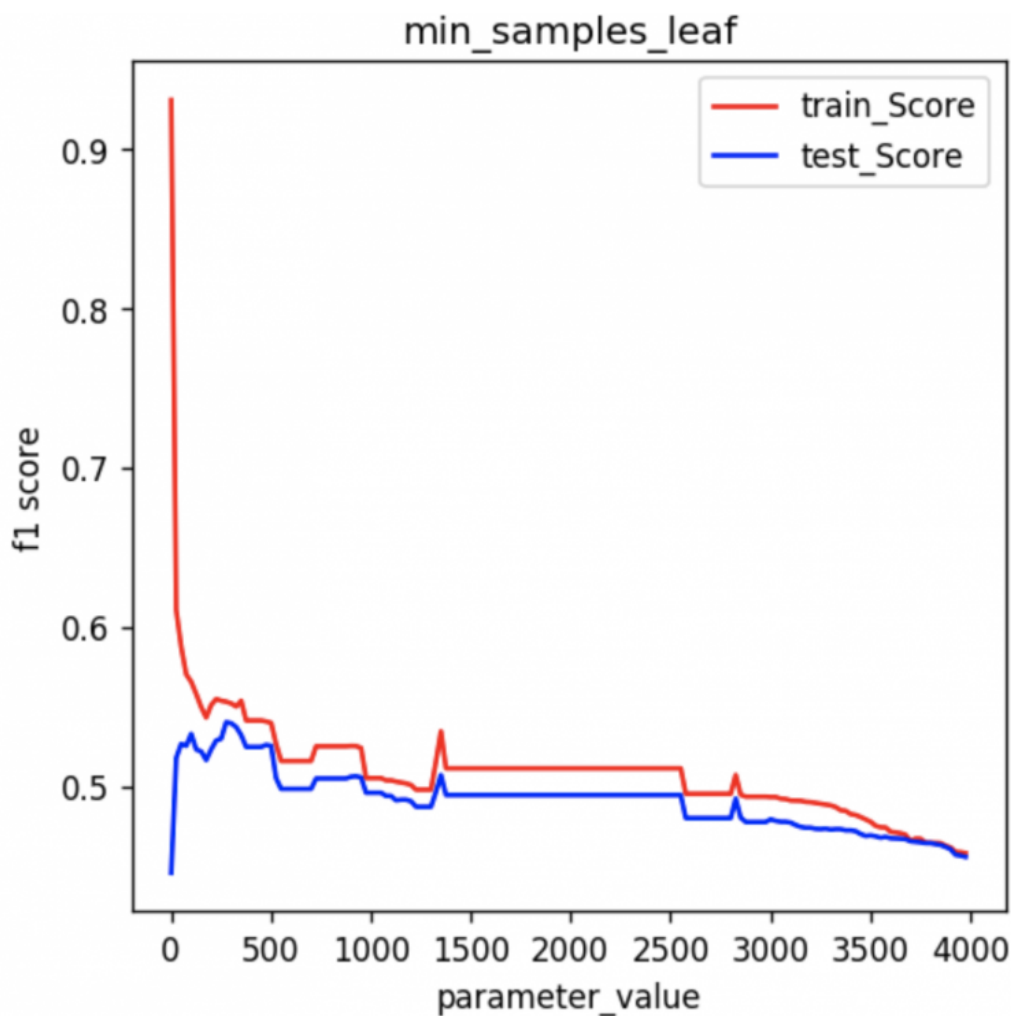
The tree on the left represents an unconstrained tree. Here, the nodes marked with green color satisfy the condition as they have a minimum of 5 samples. Hence, they will be treated as the leaf or terminal nodes.

However, the red node has only 3 samples and hence it will not be considered as the leaf node. Its parent node will become the leaf node. That's why the tree on the right represents the results when we set the minimum samples for the terminal node as 5.

So, we have controlled the growth of the tree by setting a minimum sample criterion for terminal nodes. As you would have guessed, similar to the two hyperparameters mentioned above, this hyperparameter also helps prevent overfitting as the parameter value increases.

If we plot the performance/parameter value plot as before:





(<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-15.41.51.png>).

We can clearly see that the Random Forest model is overfitting when the parameter value is very low (when parameter value < 100), but the model performance quickly rises up and rectifies the issue of overfitting (100 < parameter value < 400). But when we keep on increasing the value of the parameter (> 500), the model slowly drifts towards the realm of underfitting.

So far, we have looked at the hyperparameters that are also covered in [Decision Trees](#)

([https://courses.analyticsvidhya.com/courses/getting-started-with-decision-trees?](https://courses.analyticsvidhya.com/courses/getting-started-with-decision-trees?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning)

[utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning](https://courses.analyticsvidhya.com/courses/getting-started-with-decision-trees?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning)). Let's now look at the hyperparameters that are exclusive to Random Forest. Since Random Forest is a collection of decision trees, let's begin with the number of estimators.

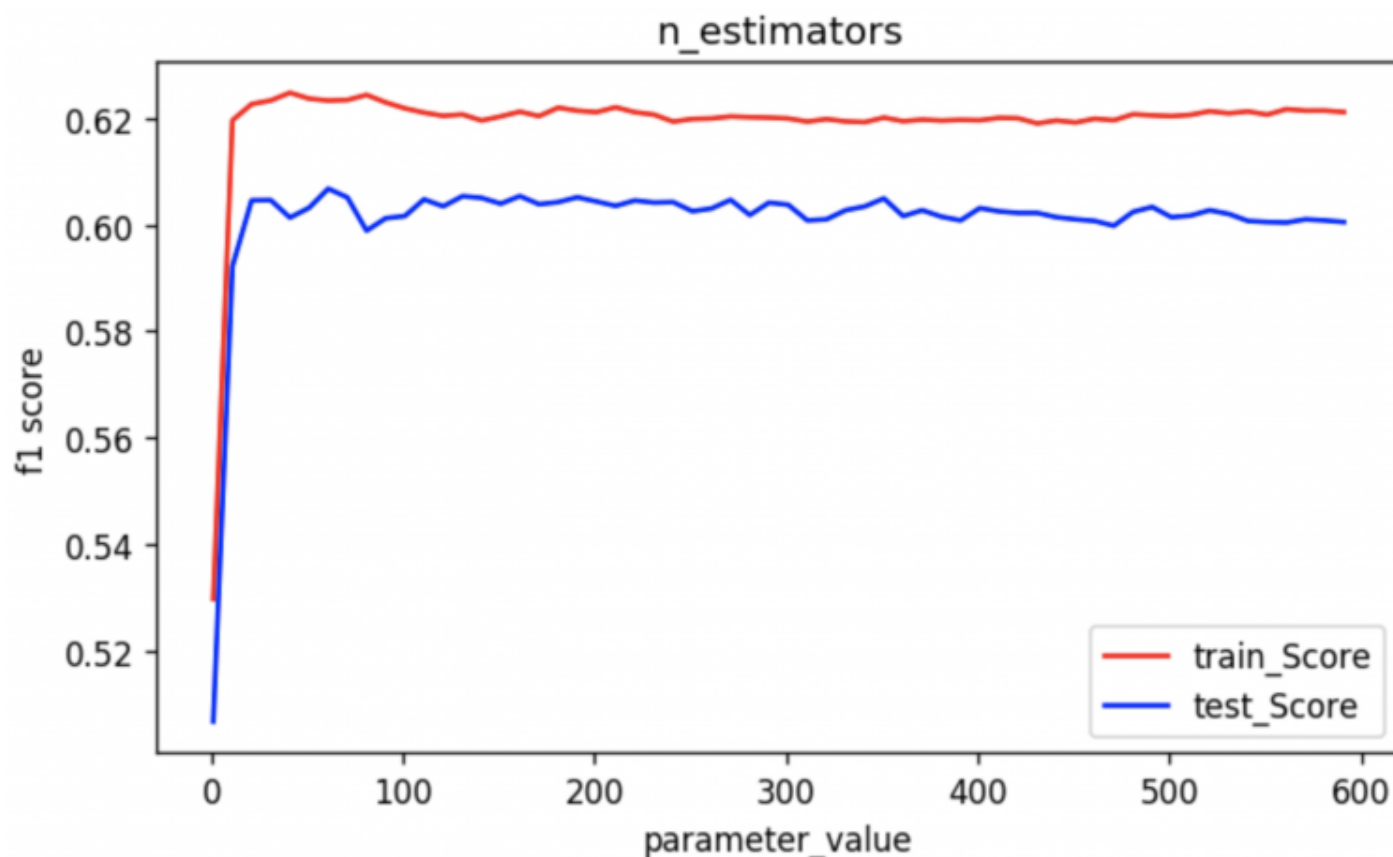
Random Forest Hyperparameter #5: n_estimators



We know that a Random Forest algorithm is nothing but a grouping of trees. But how many trees should we consider? That's a common question fresher data scientists ask. And it's a valid one!

We might say that more trees should be able to produce a more generalized result, right? But by choosing more number of trees, the time complexity of the Random Forest model also increases.

In this graph, we can clearly see that the performance of the model sharply increases and then stagnates at a certain level:

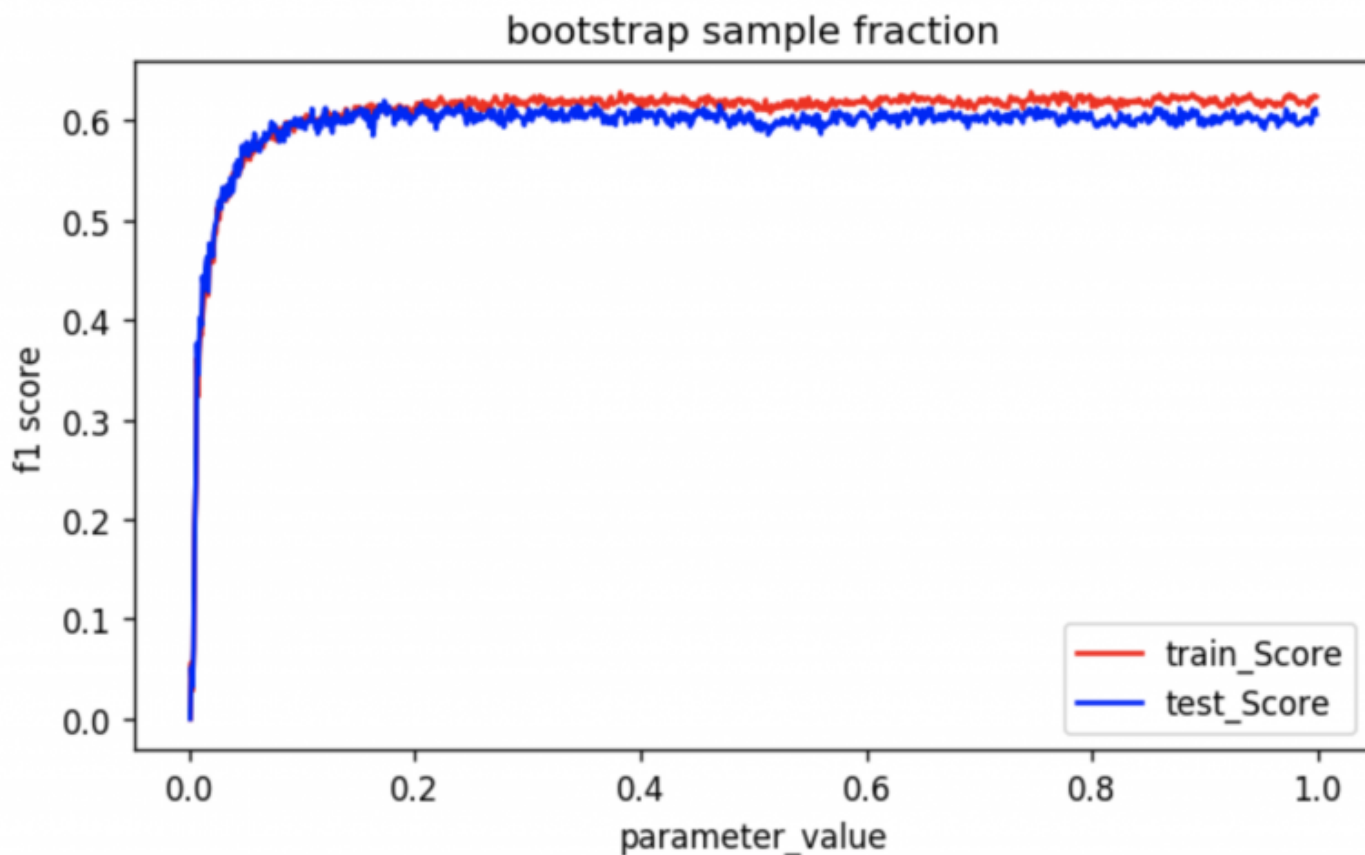


(<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-15.46.11.png>).

This means that choosing a large number of estimators in a random forest model is not the best idea. Although it will not degrade the model, it can save you the computational complexity and prevent the use of a fire extinguisher on your CPU!

Random Forest Hyperparameter #6: max_samples

The *max_samples* hyperparameter determines what fraction of the original dataset is given to any individual tree. You might be thinking that more data is always better. Let's try to see if that makes sense.



(<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-16.04.32.png>).

We can see that the performance of the model rises sharply and then saturates fairly quickly. Can you figure out what the key takeaway from this visualization is?

It is not necessary to give each decision tree of the Random Forest the full data. If you would notice, the model performance reaches its max when the data provided is less than 0.2 fraction of the original dataset. That's quite astonishing!

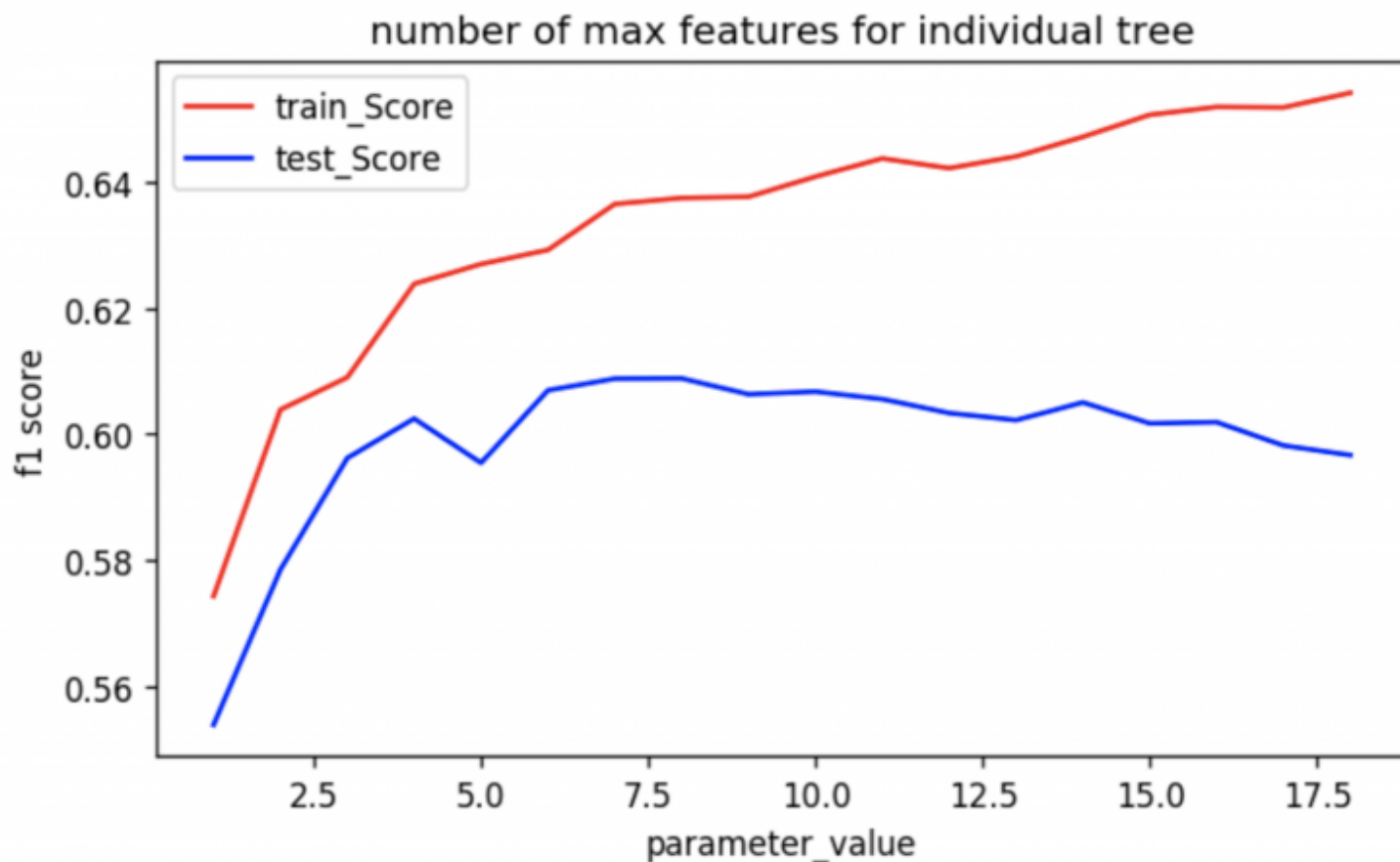
Although this fraction will differ from dataset to dataset, we can allocate a lesser fraction of bootstrapped data to each decision tree. As a result, the training time of the Random Forest model is reduced drastically.

Random Forest Hyperparameter #7: max_features

Finally, we will observe the effect of the *max_features* hyperparameter. This resembles the number of maximum features provided to each tree in a random forest.



We know that random forest chooses some random samples from the features to find the best split. Let's see how varying this parameter can affect our random forest model's performance.



(<https://cdn.analyticsvidhya.com/wp-content/uploads/2020/03/Screenshot-2020-03-04-at-16.06.57.png>).

We can see that the performance of the model initially increases as the number of *max_feature* increases. But, after a certain point, the *train_score* keeps on increasing. But the *test_score* saturates and even starts decreasing towards the end, which clearly means that the model starts to overfit.

Ideally, the overall performance of the model is the highest close to 6 value of the max features. It is a good convention to consider the default value of this parameter, which is set to square root of the number of features present in the dataset. The ideal number of max_features generally tend to lie close to this value.

End Notes

With this, we conclude our discussion on how to tune the various hyperparameters of a Random Forest model. I covered the 7 key hyperparameters here and you can explore these plus the other ones on your own. That's the best way to learn a concept and ingrain it.

Next, you should check out the comprehensive and popular [Applied Machine Learning course](https://courses.analyticsvidhya.com/courses/applied-machine-learning-beginner-to-professional?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning) (https://courses.analyticsvidhya.com/courses/applied-machine-learning-beginner-to-professional?utm_source=blog&utm_medium=beginners-guide-random-forest-hyperparameter-tuning) as the logical step in your machine learning journey! You can also read this article on Analytics Vidhya's Android APP



[. \(https://play.google.com/store/apps/details?](https://play.google.com/store/apps/details?id=com.analyticsvidhya.android&utm_source=blog_article&utm_campaign=blog&pcampaignid=MKT-Other-global-all-co-prtnr-py-PartBadge-Mar2515-1)

[id=com.analyticsvidhya.android&utm_source=blog_article&utm_campaign=blog&pcampaignid=MKT-Other-global-all-co-prtnr-py-PartBadge-Mar2515-1\)](https://play.google.com/store/apps/details?id=com.analyticsvidhya.android&utm_source=blog_article&utm_campaign=blog&pcampaignid=MKT-Other-global-all-co-prtnr-py-PartBadge-Mar2515-1)

Share this:

 (<https://www.analyticsvidhya.com/blog/2020/03/beginners-guide-random-forest-hyperparameter-tuning/?share=linkedin&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2020/03/beginners-guide-random-forest-hyperparameter-tuning/?share=facebook&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2020/03/beginners-guide-random-forest-hyperparameter-tuning/?share=twitter&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2020/03/beginners-guide-random-forest-hyperparameter-tuning/?share=pocket&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2020/03/beginners-guide-random-forest-hyperparameter-tuning/?share=reddit&nb=1>)

Related Articles



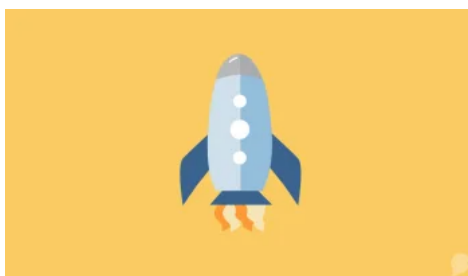
(<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/>).

Tuning the parameters of your Random Forest model

(<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/>)

June 9, 2015

In "Algorithm"



(<https://www.analyticsvidhya.com/blog/2020/02/4-boosting-algorithms-machine-learning/>).

4 Boosting Algorithms You Should Know - GBM, XGBoost, LightGBM & CatBoost

(<https://www.analyticsvidhya.com/blog/2020/02/4-boosting-algorithms-machine-learning/>)



(<https://www.analyticsvidhya.com/blog/2016/04/tree-based-algorithms-complete-tutorial-scratch-in-python/>).

Tree Based Algorithms: A Complete Tutorial from Scratch (in R & Python)

(<https://www.analyticsvidhya.com/blog/2016/04/tree-based-algorithms-complete-tutorial-scratch-in-python/>)

April 12, 2016

February 13, 2020
In "Intermediate"

In "Algorithm"

All-In-One HR Software

Intuitive HR Apps. Simplify HR Management, Reduce Busywork, Save Time & Money. Free Trial. Zenefits

TAGS : [BIAS VARIANCE TRADEOFF \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/BIAS-VARIANCE-TRADEOFF/\)](https://www.analyticsvidhya.com/blog/tag/bias-variance-tradeoff/), [CLASSIFICATION \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/CLASSIFICATION/\)](https://www.analyticsvidhya.com/blog/tag/classification/), [DECISION TREE \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/DECISION-TREE/\)](https://www.analyticsvidhya.com/blog/tag/decision-tree/), [HYPER PARAMETER TUNING \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/HYPER-PARAMETER-TUNING/\)](https://www.analyticsvidhya.com/blog/tag/hyper-parameter-tuning/), [RANDOM FOREST \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/TAG/RANDOM-FOREST/\)](https://www.analyticsvidhya.com/blog/tag/random-forest/)

NEXT ARTICLE

Top Highlights from TensorFlow Dev Summit 2020!

(<https://www.analyticsvidhya.com/blog/2020/03/highlights-tensorflow-dev-summit-2020/>)

...

PREVIOUS ARTICLE

What are Lambda Functions? A Quick Guide to Lambda Functions in Python

(<https://www.analyticsvidhya.com/blog/2020/03/what-are-lambda-functions-in-python/>)



(<https://www.analyticsvidhya.com/blog/author/sharoon-saxena/>)

[Sharoon Saxena \(Https://Www.Analyticsvidhya.Com/Blog/Author/Sharoon-Saxena/\)](https://www.analyticsvidhya.com/blog/author/sharoon-saxena/)



Passionate about learning new things everyday, well versed with Machine Learning and Data Science and an Avid Reader. Setting sights on Reinforcement Learning and Game Theory, I could see Artificial General Intelligence on the Horizon.

in_(<https://www.linkedin.com/in/sharoon-saxena-0539a0126/>).

 (<https://github.com/SharoonSaxena>).

LEAVE A REPLY

Your email address will not be published.

Comment

Name (required)

Email (required)

Website

☐ Notify me of new posts by email.

SUBMIT COMMENT





























(<https://www.analyticsvidhya.com/>)

Download App



(<https://play.google.com/store/apps/details?id=com.analyticsvidhya.android>)



(<https://apps.apple.com/us/app/analytics-vidhya/id1470025572>)

Analytics Vidhya

About Us (<https://www.analyticsvidhya.com/about-me/>)

Our Team (<https://www.analyticsvidhya.com/about-me/team/>)

Careers (<https://www.analyticsvidhya.com/about-me/career-analytics-vidhya/>)

Contact us (<https://www.analyticsvidhya.com/contact/>)

Data Science

Blog (<https://www.analyticsvidhya.com/blog/>)

Hackathon (<https://datahack.analyticsvidhya.com/>)

Discussions (<https://discuss.analyticsvidhya.com/>)

Apply Jobs (<https://www.analyticsvidhya.com/jobs/>)

Companies

Post Jobs (<https://www.analyticsvidhya.com/corporate/>)

Trainings (<https://courses.analyticsvidhya.com/>)

Hiring Hackathons (<https://datahack.analyticsvidhya.com/>)

Advertising (<https://www.analyticsvidhya.com/contact/>)

Visit us

in



(<https://www.facebook.com/analyticsvidhya/>)



(<https://www.linkedin.com/company/analytics-vidhya/>)

© Copyright 2013-2020 Analytics Vidhya

[Privacy Policy](#) [Terms of Use](#) [Refund Policy](#)

x

-

(<http://play.google.com/store/apps/details?id=com.analyticsvidhya.android>)

...

