

[f \(https://www.facebook.com/AnalyticsVidhya\)](https://www.facebook.com/AnalyticsVidhya)[t \(https://twitter.com/analyticsvidhya\)](https://twitter.com/analyticsvidhya)[g+ \(https://plus.google.com/AnalyticsVidhya\)](https://plus.google.com/AnalyticsVidhya)[in \(https://www.linkedin.com/groups/Analytics-Vidhya-Learn-everything-about-5057165\)](https://www.linkedin.com/groups/Analytics-Vidhya-Learn-everything-about-5057165)[Home \(https://www.analyticsvidhya.com/\)](https://www.analyticsvidhya.com/)[Blog \(https://www.analyticsvidhya.com/blog/\)](https://www.analyticsvidhya.com/blog/)[Jobs \(https://www.analyticsvidhya.com/jobs/\)](https://www.analyticsvidhya.com/jobs/)[Trainings \(https://www.analyticsvidhya.com/trainings/\)](https://www.analyticsvidhya.com/trainings/)[Learning Paths \(https://www.analyticsvidhya.com/learning-paths-data-science-business-analytics-business-analytics/\)](https://www.analyticsvidhya.com/learning-paths-data-science-business-analytics-business-analytics/)[DataHack \(https://datahack.analyticsvidhya.com\)](https://datahack.analyticsvidhya.com)<https://www.analyticsvidhya.com>[HOME \(HTTPS://WWW.ANALYTICSVIDHYA.COM/\)](https://www.analyticsvidhya.com/)[BLOG \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/\)](https://www.analyticsvidhya.com/blog/)[TRAININGS \(HTTPS://WWW.ANALYTICSVIDHYA.COM/TRAININGS/\)](https://www.analyticsvidhya.com/trainings/)[DISCUSS \(HTTPS://DISCUSS.ANALYTICSVIDHYA.COM/\)](https://discuss.analyticsvidhya.com/)[LEARNING PATHS \(HTTPS://WWW.ANALYTICSVIDHYA.COM/LEARNING-PATHS-DATA-SCIENCE-BUSINESS-ANALYTICS-BUSINESS-ANALYTICS/\)](https://www.analyticsvidhya.com/learning-paths-data-science-business-analytics-business-analytics/)[DATAHACK \(HTTPS://DATAHACK.ANALYTICSVIDHYA.COM\)](https://datahack.analyticsvidhya.com) [STORIES \(HTTPS://WWW.ANALYTICSVIDHYA.COM/STORIES/\)](#)[WRITE FOR US \(HTTP://WWW.ANALYTICSVIDHYA.COM/ABOUT-ME/WRITE/\)](http://www.analyticsvidhya.com/about-me/write/)[CONTACT US \(HTTPS://WWW.ANALYTICSVIDHYA.COM/CONTACT-US/\)](https://www.analyticsvidhya.com/contact-us/)

[Home \(https://www.analyticsvidhya.com/\)](https://www.analyticsvidhya.com/) > [Business Analytics \(https://www.analyticsvidhya.com/blog/category/business-analytics/\)](https://www.analyticsvidhya.com/blog/category/business-analytics/) > [Tuning the parameters of your Random Forest model \(https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/\)](https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/)

# Tuning the parameters of your Random Forest model

[BUSINESS ANALYTICS \(HTTPS://WWW.ANALYTICSVIDHYA.COM/BLOG/CATEGORY/BUSINESS-ANALYTICS/\)](https://www.analyticsvidhya.com/blog/category/business-analytics/)[SHARE f \(http://www.facebook.com/sharer.php?u=https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/\)](http://www.facebook.com/sharer.php?u=https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/)

model/&t=Tuning%20the%20parameters%20of%20your%20Random%20Forest%20model)  (<https://twitter.com/home?status=Tuning%20the%20parameters%20of%20your%20Random%20Forest%20model-https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/>)  (<https://plus.google.com/share?url=https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/>)  (<http://pinterest.com/pin/create/button/?url=https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/&media=https://www.analyticsvidhya.com/wp-content/uploads/2015/06/Tuning-the-parameters-of-your-Random-Forest-model.jpg&description=Tuning%20the%20parameters%20of%20your%20Random%20Forest%20model>)

## Why to tune Machine Learning Algorithms?

A month back, I participated in a Kaggle competition (<https://www.analyticsvidhya.com/blog/2015/06/start-journey-kaggle/>) called TFI. I started my first submission at 50th percentile. Having worked relentlessly on feature engineering for 3 weeks, I managed to reach 20th percentile. To my surprise, right after tuning the parameters of the machine learning algorithm I was using, I was able to breach top 10th percentile.

This is how important tuning these machine learning algorithms are. Random Forest is the easiest machine learning tool used in the industry. In our previous articles, we have introduced Random Forest (<https://www.analyticsvidhya.com/blog/2014/06/introduction-random-forest-simplified/>) and compared it against a CART algorithm (<https://www.analyticsvidhya.com/blog/2014/06/comparing-cart-random-forest-1/>). Machine Learning tools are known for their performance.

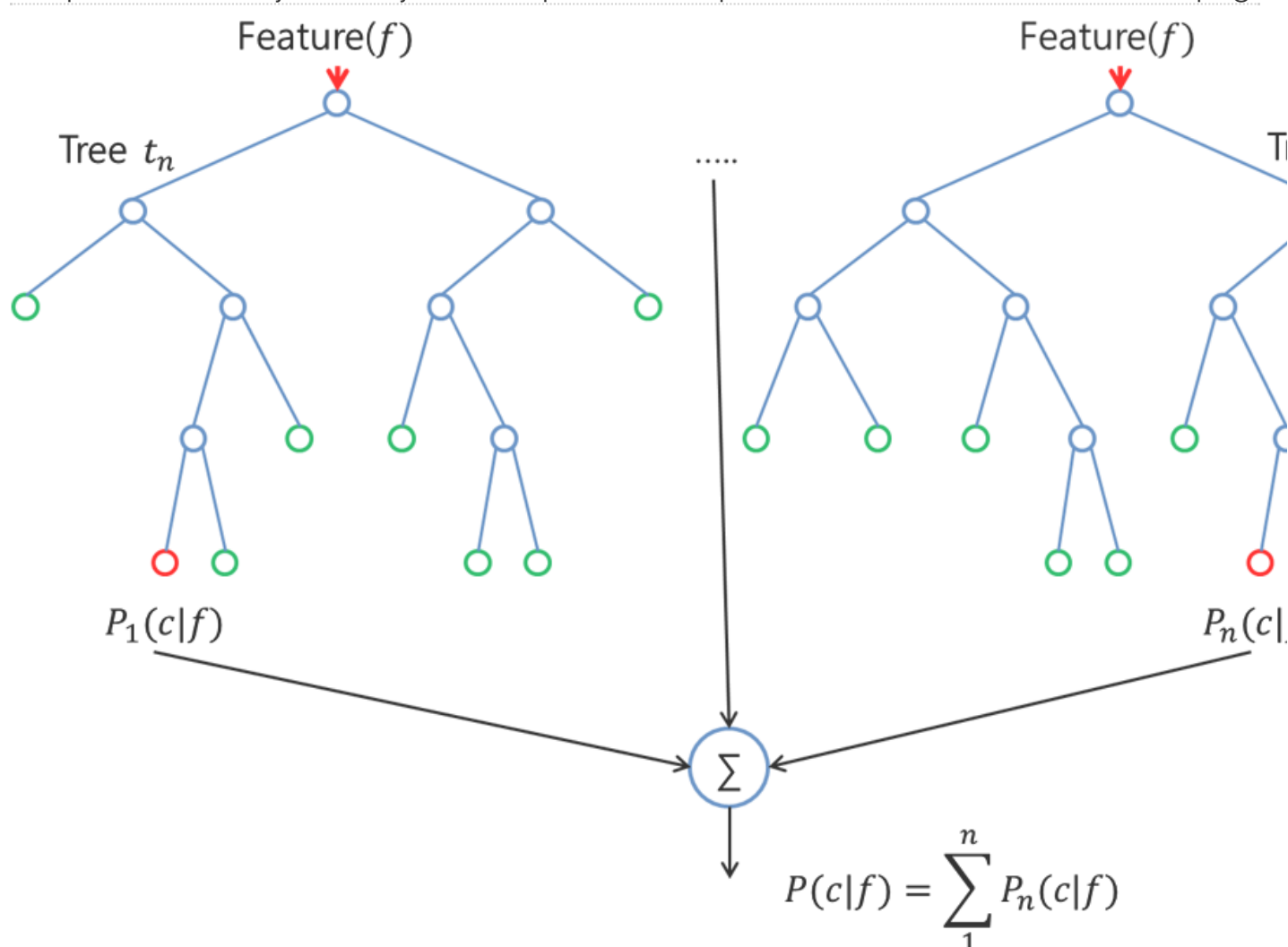


(<https://www.analyticsvidhya.com/wp-content/uploads/2015/06/Tuning-the-parameters-of-your-Random-Forest-model.jpg>)

# What is a Random Forest?

Random forest is an ensemble tool which takes a subset of observations and a subset of features to build a decision tree. It builds multiple such decision trees and amalgamates them together to give a more accurate and stable prediction. This is a direct consequence of the fact that by maximizing the accuracy of a panel of independent judges, we get the final prediction better than the best judge.

(<https://www.analyticsvidhya.com/wp-content/uploads/2015/06/random-forest4.png>)



(<https://www.analyticsvidhya.com/wp-content/uploads/2015/06/random-forest7.png>)

We generally see a random forest as a black box which takes in input and gives out prediction without worrying too much about what calculations are going on the back end. This black box has a few levers we can play with. Each of these levers has some effect on either the prediction accuracy or the model's interpretability.

of the model or the resource – time balance. In this article we will talk more about these can tune, while building a random forest model.

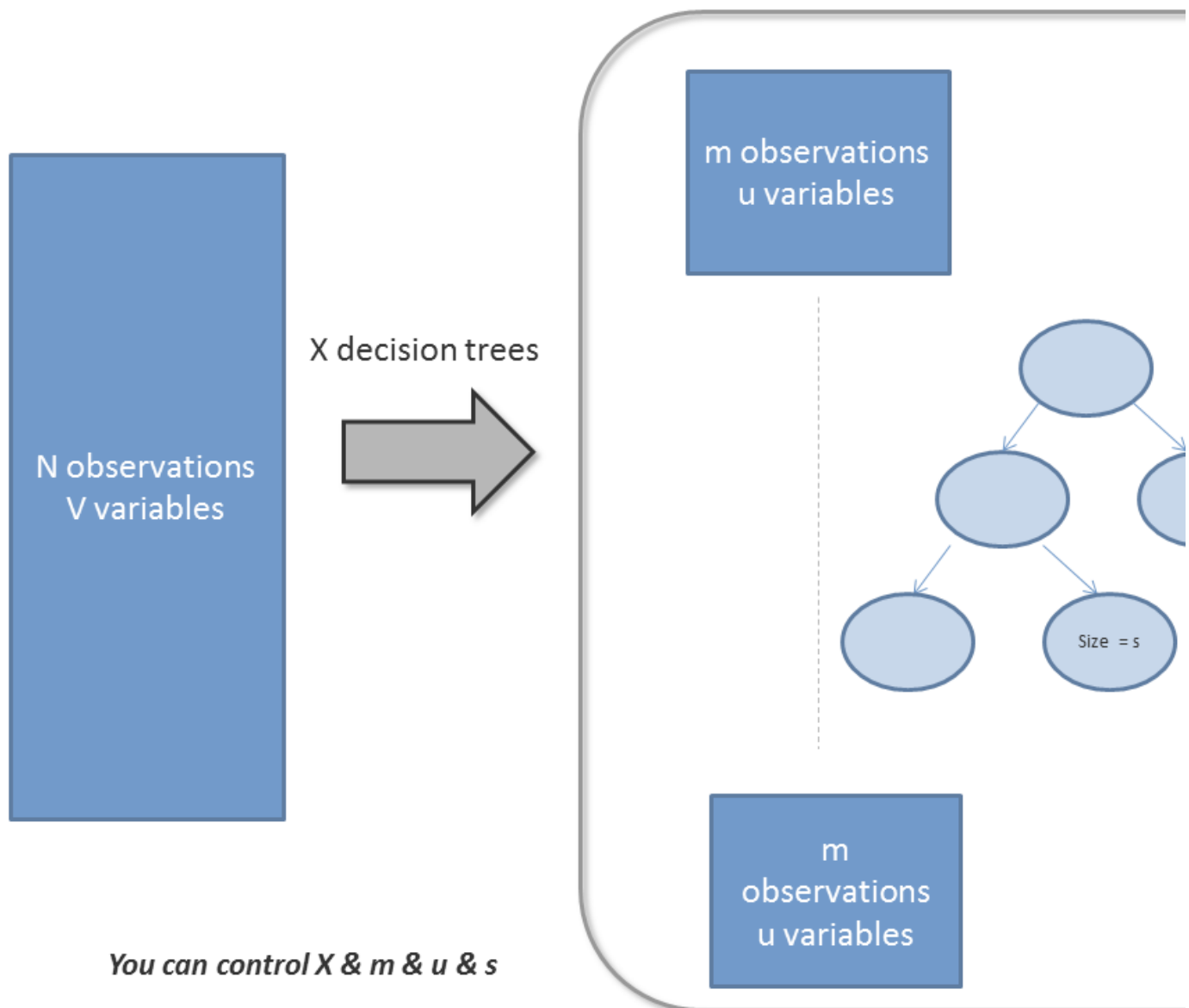
## Parameters / levers to tune Random Forests

Parameters in random forest are either to increase the predictive power of the model or easier to train the model. Following are the parameters we will be talking about in more detail (that I am using Python conventional nomenclatures for these parameters) :



(<https://www.analyticsvidhya.com/wp-content/uploads/2015/06/knobs.jpg>)

# 1. Features which make predictions of the model better



(<https://www.analyticsvidhya.com/wp-content/uploads/2015/06/RF.png>)

There are primarily 3 features which can be tuned to improve the predictive power of the model.

## 1.a. max\_features:

These are the maximum number of features Random Forest is allowed to try in individual trees. There are multiple options available in Python to assign maximum features. Here are a few of them:

1. *Auto/None* : This will simply take all the features which make sense in every tree. Here we do not put any restrictions on the individual tree.
2. *sqrt* : This option will take square root of the total number of features in individual run. For example, if the total number of variables are 100, we can only take 10 of them in individual tree.  $\log_2(100) \approx 6.64$ , so we take the integer part, which is 6, and then take the square root of 6, which is approximately 2.45, so we take the integer part, which is 2.

similar type of option for max\_features.

3. 0.2 : This option allows the random forest to take 20% of variables in individual run. We can use a value in a format "0.x" where we want x% of features to be considered.

### **How does "max\_features" impact performance and speed?**

Increasing max\_features generally improves the performance of the model as at each node we have a higher number of options to be considered. However, this is not necessarily true as it decreases the diversity of individual tree which is the USP of random forest. But, it also decreases the speed of algorithm by increasing the max\_features. Hence, you need to strike a balance and choose the optimal max\_features.

### **1.b. n\_estimators :**

This is the number of trees you want to build before taking the maximum voting or average for predictions. Higher number of trees give you better performance but makes your code slower. You should choose as high a value as your processor can handle because this makes your model stronger and more stable.

### **1.c. min\_sample\_leaf :**

If you have built a decision tree before, you can appreciate the importance of minimum sample size. Leaf is the end node of a decision tree. A smaller leaf makes the model more prone to overfitting due to noise in train data. Generally I prefer a minimum leaf size of more than 50. However, you should try multiple leaf sizes to find the most optimum for your use case.

## **2. Features which will make the model training easier**

There are a few attributes which have a direct impact on model training speed. Following are the parameters which you can tune for model speed :

### **2.a. n\_jobs :**

This parameter tells the engine how many processors it is allowed to use. A value of "-1" means no restriction whereas a value of "1" means it can only use one processor. Here is an experiment you can do with Python to check this metric :

```
%timeit
```

```
model = RandomForestRegressor(n_estimators = 100, oob_score = TRUE, n_jobs = 1, random_
```

```
model.fit(X,y)
```

Output ——— 1 loop best of 3 : 1.7 sec per loop

```
%timeit
```

```
model = RandomForestRegressor(n_estimators = 100, oob_score = TRUE, n_jobs = -1, random_
```

```
model.fit(X,y)
```

Output ——— 1 loop best of 3 : 1.1 sec per loop

"%timeit" is an awesome function which runs a function multiple times and gives the fastest time. This comes out very handy while scaling up a particular function from prototype to fir

## 2.b. random\_state :

This parameter makes a solution easy to replicate. A definite value of random\_state produces the same results if given with the same parameters and training data. I have personally created an ensemble with multiple models of different random states and all optimum parameters perform better than individual random state.

## 2.c. oob\_score :

This is a random forest cross validation method. It is very similar to leave one out validation however, this is so much faster. This method simply tags every observation used in different folds. Then it finds out a maximum vote score for every observation based on only trees which did not use this particular observation to train itself.

Here is a single example of using all these parameters in a single function :

```
model = RandomForestRegressor(n_estimator = 100, oob_score = TRUE, n_jobs = -1, random_state = 50, max_features = "auto", min_samples_leaf = 5)

model.fit(X,y)
```

## Learning through a case study

We have referred to Titanic case study in many of our previous articles. Let's try the same again. The objective of this case here will be to get a feel of random forest parameter tuning and getting the right features. Try following code to build a basic model :

```
from sklearn.ensemble import RandomForestRegressor
```

```
from sklearn.metrics import roc_auc_score
```

```
import pandas as pd
```

```
x = pd.read_csv("train.csv")
```

```
y = x.pop("Survived")
```



```
model = RandomForestRegressor(n_estimator = 100 , oob_score = TRUE, random_state = 1)

model.fit(x(numeric_variable,y))

print "AUC - ROC : ", roc_auc_score(y,model.oob_prediction)
```

AUC – ROC : 0.7386

This is a very simplistic model with no parameter tuning. Now let's do some parameter tuning. As we have discussed before, we have 6 key parameters to tune. We have some grid search algorithm in Python, which can tune all parameters automatically. But here let's get our hands dirty and understand the mechanism better. Following code will help you tune the model for different parameters.

***Exercise : Try running the following code and find the optimal leaf size in the comment below.***

```
sample_leaf_options = [1,5,10,50,100,200,500]

for leaf_size in sample_leaf_options :

    model = RandomForestRegressor(n_estimator = 200, oob_score = TRUE, n_jobs = -1,
    leaf_size = 50, max_features = "auto", min_samples_leaf = 10,
    min_samples_split = 10, random_state = 1,
    leaf_size)

    model.fit(x(numeric_variable,y))

    print "AUC - ROC : ", roc_auc_score(y,model.oob_prediction)
```

## End Notes


Machine learning tools like random forest, SVM, neural networks etc. are all used for high performance. They do give high performance, but users generally don't understand how they work. Not knowing the statistical details of the model is not a concern however not knowing how the model can be tuned well to clone the training data restricts the user to use the algorithm's full potential. In some of the future articles we will take up tuning of other machine learning algorithms like SVM, GBM and neural networks.

Have you used random forest before? What parameters did you tune? How did tuning them impact the performance of the model? Did you see any significant benefits by doing the same? Let us know your thoughts about this guide in the comments section below.

**If you like what you just read & want to continue your analytics learning, subscribe to our emails (<http://feedburner.google.com/fb/a/mailverify?uri=analyticsvidhya>), follow us on twitter (<http://twitter.com/analyticsvidhya>) or like our facebook page (<http://facebook.com/analyticsvidhya>).**

---

### Share this:

 (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?share=linkedin&nb=1>) 432

 (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?share=facebook&nb=1>) 174

 (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?share=google-plus-1&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?share=twitter&nb=1>)

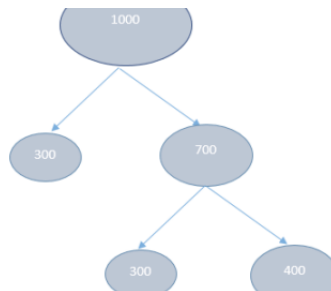
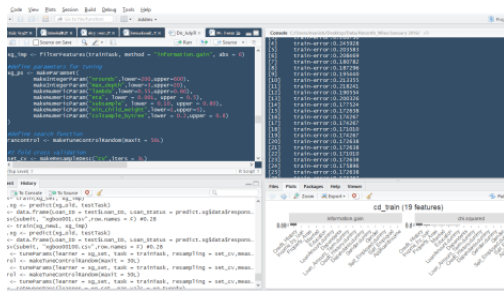
 (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?share=pocket&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?share=reddit&nb=1>)

---

## RELATED

---



(<https://www.analyticsvidhya.com/blog/2016/08/practicing-machine-learning-techniques-in-r-with-mlr-package/>)

**Practicing Machine Learning Techniques in R with MLR Package**  
(<https://www.analyticsvidhya.com/blog/2016/08/practicing-machine-learning-techniques-in-r-with-mlr-package/>)

August 8, 2016

In "Machine Learning"

(<https://www.analyticsvidhya.com/blog/2016/12/detailed-solutions-for-skilltest-tree-based-algorithms/>)

**45 questions to test Data Scientists on Tree Based Algorithms (Decision tree, Random Forests, XGBoost)**

(<https://www.analyticsvidhya.com/blog/2016/12/detailed-solutions-for-skilltest-tree-based-algorithms/>)

December 5, 2016

In "Machine Learning"

(<https://www.analyticsvidhya.com/blog/2016/05/h2o-data-table-models-large-data-sets/>)

**Use H2O and data.table to build models on large data sets**

(<https://www.analyticsvidhya.com/blog/2016/05/h2o-data-table-models-large-data-sets/>)

May 12, 2016

In "Machine Learning"

TAGS: MAX\_FEATURES (<https://www.analyticsvidhya.com/blog/tag/max-features/>), MIN\_SAMPLE\_LEAF

(<https://www.analyticsvidhya.com/blog/tag/min-sample-leaf/>), N\_ESTIMATOR (<https://www.analyticsvidhya.com/blog/tag/n-estimator/>)

N\_JOBS (<https://www.analyticsvidhya.com/blog/tag/n-jobs/>), PARAMETER\_TUNING (<https://www.analyticsvidhya.com/blog/tag/parameter-tuning/>), PYTHON (<https://www.analyticsvidhya.com/blog/tag/python/>), RANDOM FOREST (<https://www.analyticsvidhya.com/blog/tag/random-forest/>)

FOREST/)



**Previous Article**

**Cheat Sheet for Exploratory Data Analysis in Python**

(<https://www.analyticsvidhya.com/blog/2015/06/infographic-cheat-sheet-data-exploration-python/>)

**Next Article**

**In Conversation with Mr. Stefan Groschupf, Founder and CEO, Datameer**  
(<https://www.analyticsvidhya.com/blog/2016/05/in-conversation-with-mr-stefan-groschupf-ceo-datameer/>)



(<https://www.analyticsvidhya.com/blog/author/tavish1/>)

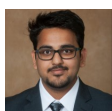
Author

**Tavish Srivastava**

(<https://www.analyticsvidhya.com/blog/author/tavish1/>)

I am Tavish Srivastava, a post graduate from IIT Madras in Mechanical Engineering. I have more than two years of work experience in Analytics. My experience ranges from hands on analytics in a developing country like India to convince banking partners with analytical solution in matured market like US. For last two and a half years I have contributed to various sales strategies, marketing strategies and Recruitment strategies in both Insurance and Banking industry.

## 5 COMMENTS



**Aayush Agrawal says:**

REPLY (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?replytocom=88206>)  
JUNE 10, 2015 AT 6:50 AM (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/>)

Brilliantly written article. Currently I have used all of these techniques in a Data science project I was working on and it definitely helps in improving model performance and accuracy. I also came across something else also when I was reading some articles on Random Forest, Regularization of Random Forest. The theme was to not only split data with some variables but also cross-validation. Splitting is significant enough using Statistical validation, now this is something which can help in taking Random Forest to next level, as It can help in reducing over-fitting. I tried to use caret package but I think this technique is computationally expensive so couldn't run it on my system. I would love to see an article on it to understand it's working and how its performance can be improved.



**KARTHI V says:**

REPLY (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?replytocom=88254>)  
JUNE 10, 2015 AT 4:00 PM (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/>)

Hi Tavish,

Very useful article.



**Ravi says:**

REPLY ([https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?replytocom=SEPTMBER 11, 2015 AT 6:43 PM](https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?replytocom=SEPTMBER%2015%20AT%206%3A43%20PM) (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/#comment-94729>))

I love AV and am a fan of your articles. I have heard something like Conditional Inference which are similar to Random Forests. Can you share your thoughts on Conditional Inference also? How does it work & its tuning parameters, when does it outcast Random Forests?



**Josh says:**

REPLY ([https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?replytocom=SEPTMBER 17, 2015 AT 2:39 AM](https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?replytocom=SEPTMBER%2017%20AT%202%3A39%20AM) (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/#comment-95226>))

Great article! I would love to see something similar regarding parameter tuning for the r package.



**John S says:**

REPLY ([https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?replytocom=OCTOBER 12, 2016 AT 12:15 AM](https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/?replytocom=OCTOBER%2012%20AT%2012%3A15%20AM) (<https://www.analyticsvidhya.com/blog/2015/06/tuning-random-forest-model/#comment-117085>))

This was a very nice article. I would still be interested to know if there is a minimum number of trees that can be calculated to reduce computational cost?

## LEAVE A REPLY

Connect with:



([https://www.analyticsvidhya.com/wp-login.php?](https://www.analyticsvidhya.com/wp-login.php?action=wordpress_social_authenticate&mode=login&provider=Facebook&redirect_to=https%3A%2F%2Fwww.analyticsvidhya.com%2Fblog%2F2015%2F06%2Ftuning-random-forest-model%2F)

[action=wordpress\\_social\\_authenticate&mode=login&provider=Facebook&redirect\\_to=https%3A%2F%2Fwww.analyticsvidhya.com%2Fblog%2F2015%2F06%2Ftuning-random-forest-model%2F](https://www.analyticsvidhya.com/wp-login.php?action=wordpress_social_authenticate&mode=login&provider=Facebook&redirect_to=https%3A%2F%2Fwww.analyticsvidhya.com%2Fblog%2F2015%2F06%2Ftuning-random-forest-model%2F))

Your email address will not be published.

Comment

Name (required)

Email (required)

Website

SUBM

## ABOUT US

For those of you, who are wondering what is "Analytics Vidhya", "Analytics" can be defined as the science of extracting insights from raw data. The spectrum of analytics starts from capturing data and evolves into using insights / trends from this data to make informed decisions. [Read More](http://www.analyticsvidhya.com/about-me/) (<http://www.analyticsvidhya.com/about-me/>)

## STAY CONNECTED

## LATEST POSTS



(<https://www.analyticsvidhya.com/blog/plan-2017-beginners-data-science/>)

**Infographic - Learning Plan 2  
beginners in data science**  
(<https://www.analyticsvidhya.com/blog/2017/01/learning-plan-2017-b-data-science/>)

KUNAL JAIN , JANUARY 28, 2017



(<http://www.twitter.com/andrewswahy>) (<https://www.facebook.com/AndrewSwahy>)



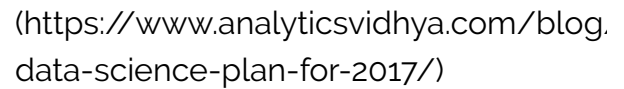
(<https://www.facebook.com/ArcticSwindly>)



(<https://plus.google.com/+Atlasyifoodbyme>) (<https://www.google.com/fb/a/mailverifier>)



```
uri=analyticsvidhya)
```



KUNAL JAIN , JANUARY 28, 2017



KUNAL JAIN , JANUARY 28, 2017



YOGESH KULKARNI . JANUARY 27, 20'