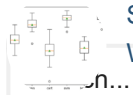Save 30% or More With Our Ebook Bundles

# Machine Learning Mastery
## Making Developers Awesome at Machine Learning
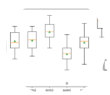
How to Develop Multi-Output Regression
Models with Python

Click to Take the FREE Ensemble Learning Crash-Course

Stacking Ensemble Machine Learning
With Python

# How to Develop a Random Forest Ensemble in Python

How to Develop Super Learner
Ensembles in Python

by **Jason Brownlee** on April 20, 2020 in **Ensemble Learning**

How to Develop Voting Ensembles With
Python

Share        Share

Last Updated on April 27, 2021

One-vs-Rest and One-vs-One for Multi-
Class Classification

Random forest is an ensemble machine learning algorithm.

It is perhaps the most popular and widely used machine learning algorithm given its good or excellent performance across a wide range of classification and regression predictive modeling problems.

**Loving the Tutorials?**

It is also easy to use given that it has few key hyperparameters and sensible heuristics for configuring these hyperparameters.

The Ensemble Learning With Python EBook
is where you'll find the *Really Good* stuff.

In this tutorial                              to develop a random forest ensemble for classification and
regression..

    >> SEE WHAT'S INSIDE

After completing this tutorial, you will know:

- Random forest ensemble is an ensemble of d
- How to use the random forest ensemble for cl
- How to explore the effect of random forest mo

**Kick-start your project** with my new book Ensem                                            ep-
by-step tutorials and the *Python source code* files

Let's get started.

- **Update Aug/2020**: Added a common question

**Start Machine Learning**    ✕

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about

services and special offers by email. For more

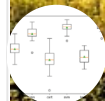information, see the Privacy Policy.

START MY EMAIL COURSE
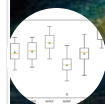
**Never miss a tutorial:**



How to Develop Multi-Output Regression Models With Python

Stacking Ensemble Machine Learning With Python

How to Develop Super Learner Ensembles in Python

How to Develop Voting Ensembles With Python

One-vs-Rest and One-vs-One for Multi-Class Classification

How to Develop a Random Forest Ensemble in Python
Photo by Sheila Sund, some rights reserved.

# Tutorial Overview

**Loving the Tutorials?**

This tutorial is divided into four parts; they are:

The Ensemble Learning With Python EBook is where you'll find the *Really Good* stuff.

1. Random Forest Algorithm
2. Random Forest Scikit-Learn API
    >> SEE WHAT'S INSIDE
    1. Random Forest for Classification
    2. Random Forest for Regression
3. Random Forest Hyperparameters
    1. Explore Number of Samples
    2. Explore Number of Features
    3. Explore Number of Trees
    4. Explore Tree Depth
4. Common Questions

# Random Forest Algorithm

Random forest is an ensemble of decision tree alg

It is an extension of bootstrap aggregation (baggin
and regression problems.

**Start Machine Learning** ✕

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about services and special offers by email. For more information, see the Privacy Policy.

**START MY EMAIL COURSE**

In bagging, a number of decision trees are created where each tree is created from a different bootstrap sample of the training dataset. A bootstrap sample is a sample of the training dataset where a sample may appear more than once in the sample, referred to as **sampling with replacement**.

Bagging is an effective ensemble algorithm as each decision tree is fit on a slightly different training dataset, and in turn, has a slightly different performance. Unlike normal decision tree models, such as classification and regression trees (CART), trees used in the ensemble are unpruned, making them slightly overfit to the training dataset. This is desirable as it helps to make each tree more different and have less correlated predictions or prediction errors.

Predictions from the trees are averaged across all decision trees resulting in better performance than any single tree in the model.

> *Each model in the ensemble is then used to generate a prediction for a new sample and these m predictions are averaged to give the forest's prediction*

— Page 199, Applied Predictive Modeling, 2013.

A prediction on a regression problem is the average of the prediction across the trees in the ensemble. A prediction on a classification problem is the majority vote for the class label across the trees in the ensemble.

- **Regression**: Prediction is the average prediction across the decision trees.
- **Classification**: Prediction is the majority vote class label predicted across the decision trees.

> *As with bagging, each tree in the forest casts a vote for the classification of a new sample, and the proportion of votes in each class across the ensemble is the predicted probability vector.*

— Page 387, Applied Predictive Modeling, 2013.

Random forest involves constructing a large number of decision trees from bootstrap samples from the training dataset, like bagging.

Unlike bagging, random forest also involves selecting a subset of input features (columns or variables) at each split point in the construction of trees. Typically, constructing a decision tree involves evaluating the value for each input variable in the data in order to select a split point. By reducing the features to a random subset that may be considered at each split point, it forces each decision tree in the ensemble to be more different.

> *Random forests provide an improvement over bagged trees by way of a small tweak that decorrelates the trees. [...] But when building these decision trees, each time a split in a tree is considered, a random sample of m predictors is chosen as split candidates from the full set of p predictors.*

— Page 320, An Introduction to Statistical Learning

The effect is that the predictions, and in turn, prediction errors, made by each tree in the ensemble are more different or less correlated. When the predictions from these less correlated trees are averaged to make a prediction, it often results in better performance than bagged decision trees.
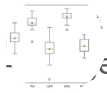
Perhaps the most important hyperparameter to tune for the random forest is the number of random features to consider at each split point.

*Random forests' tuning parameter is the number of randomly selected predictors, k, to choose from at each split, and is commonly referred to as mtry. In the regression context, Breiman (2001) recommends setting mtry to be one-third of the number of predictors.*

— Page 199, Applied Predictive Modeling, 2013.

A good heuristic for regression is to set this hyperparameter to 1/3 the number of input features.

- num_features_for_split = total_input_features / 3

*For classification problems, Breiman (2001) recommends setting mtry to the square root of the number of predictors.*

— Page 387, Applied Predictive Modeling, 2013.

A good heuristic for classification is to set this hyperparameter to the square root of the number of input features.

- num_features_for_split = sqrt(total_input_features)

Another important hyperparameter to tune is the depth of the decision trees. Deeper trees are often more overfit to the training data, but also less correlated, which in turn may improve the performance of the ensemble. Depths from 1 to 10 levels may be effective.

Finally, the number of decision trees in the ensemble can be set. Often, this is increased until no further improvement is seen.

---

## Want to Get Started W

Take my free 7-day email cras

Click to sign-up and also get a fre

Download Your

---

# Random Forest Scikit-Learn A

Random Forest ensembles can be implemented from scratch, although this can be challenging for beginners.

The scikit-learn Python machine learning library provides an implementation of Random Forest for machine learning.

It is available in modern versions of the library.

How to Develop Multi-Output Regression Models with Python

You can confirm that you are using a modern version of the library by running the following script:

```
1  # check scikit-learn version
2  import sklearn
3  print(sklearn.__version__)
```

Stacking Ensemble Machine Learning With Python

Running the script will print your version of scikit-learn.

Your version should be the same or higher. If not, you must upgrade your version of the scikit-learn

How to Develop Super Learner Ensembles in Python

```
1  0.22.1
```

How to Develop Voting Ensembles With Python

Random Forest is provided via the RandomForestRegressor and RandomForestClassifier classes.

Both models operate the same way and take the same arguments that influence how the decision trees are created.

One-vs-Rest and One-vs-One for Multi-Class Classification

Randomness is used in the construction of the model. This means that each time the algorithm is run on the same data, it will produce a slightly different model.

When using machine learning algorithms that have a stochastic learning algorithm, it is good practice to evaluate them by averaging their performance across multiple runs or repeats of cross-validation. When fitting a final model, it may be desirable to either increase the number of trees until the variance of the model is reduced across repeated evaluations, or to fit multiple final models and average their predictions.

Let's take a look at how to develop a Random Forest ensemble for both classification and regression tasks.

## Random Forest for Classification

In this section, we will look at using Random Forest

First, we can use the make_classification() function

with 1,000 examples and 20 input features.

The complete example is listed below.

```
1  # test classification dataset
2  from sklearn.datasets import make_classific
3  # define dataset
4  X, y = make_classification(n_samples=1000,
5  # summarize the dataset
6  print(X.shape, y.shape)
```

Running the example creates the dataset and summarizes the shape of the input and output components.

```
1  (1000, 20) (1000,)
```

Next, we can evaluate a random forest algorithm on this dataset.

We will evaluate the model using repeated stratified k-fold cross-validation, with three repeats and 10 folds. We will report the mean and standard deviation of the accuracy of the model across all repeats and folds.

```
1   # evaluate random forest algorithm for classification
2   from numpy import mean
3   from numpy import std
4   from sklearn.datasets import make_classification
5   from sklearn.model_selection import cross_val_score
6   from sklearn.model_selection import RepeatedStratifiedKFold
7   from sklearn.ensemble import RandomForestClassifier
8   # define dataset
9   X, y = make_classification(n_samples=1000, n_features=20, n_informative=15, n_redundant=5,
10  # define the model
11  model = RandomForestClassifier()
12  # evaluate the model
13  cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
14  n_scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1, error_score='
15  # report performance
16  print('Accuracy: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
```

Running the example reports the mean and standard deviation accuracy of the model.

**Note**: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

In this case, we can see the random forest ensemble with default hyperparameters achieves a classification accuracy of about 90.5 percent.

```
1  Accuracy: 0.905 (0.025)
```

We can also use the random forest model as a final model and make predictions for classification.

First, the random forest ensemble is fit on all available data, then the *predict()* function can be called to make predictions on new data.

The example below demonstrates this on our binary

```
1   # make predictions using random forest for
2   from sklearn.datasets import make_classifi
3   from sklearn.ensemble import RandomForestC
4   # define dataset
5   X, y = make_classification(n_samples=1000,
6   # define the model
7   model = RandomForestClassifier()
8   # fit the model on the whole dataset
9   model.fit(X, y)
10  # make a single prediction
11  row = [[-8.52381793,5.24451077,-12.1496770
12  yhat = model.predict(row)
13  print('Predicted Class: %d' % yhat[0])
```

Running the example fits the random forest ensemble model on the entire dataset and is then used to make a prediction on a new row of data, as we might when using the model in an application.

**Never miss a tutorial:**

```
1  Predicted Class: 0
```

Now that we are familiar with using random forest for classification, let's look at the API for regression.
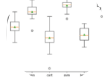
**Picked for you:**

# Random Forest for Regression

How to Develop Multi-Output Regression Models with Python

In this section, we will look at using random forests for a regression problem.

First, we can use the make_regression() function to create a synthetic regression problem with 1,000

Stacking Ensemble Machine Learning With Python

es and 20 input features.

The complete example is listed below.

How to Develop Super Learner

```
1  # test regression dataset
2  from sklearn.datasets import make_regression
3  # define dataset
4  X, y = make_regression(n_samples=1000, n_features=20, n_informative=15, noise=0.1, random_st
5  # summarize the dataset
6  print(X.shape, y.shape)
```

Python

Running the example creates the dataset and summarizes the shape of the input and output components.

One-vs-Rest and One-vs-One for Multi-Class Classification

```
1  (1000, 20) (1000,)
```

Next, we can evaluate a random forest algorithm on this dataset.

As we did with the last section, we will evaluate the model using repeated k-fold cross-validation, with three repeats and 10 folds. We will report the mean absolute error (MAE) of the model across all **Loving the Tutorials?** repeats and folds. The scikit-learn library makes the MAE negative so that it is maximized instead of minimized. This means that larger negative MAE are better and a perfect model has a MAE of 0.

The complete example is listed below.

>> SEE WHAT'S INSIDE

```
1   # evaluate random forest ensemble for regression
2   from numpy import mean
3   from numpy import std
4   from sklearn.datasets import make_regressi
5   from sklearn.model_selection import cross_
6   from sklearn.model_selection import Repeat
7   from sklearn.ensemble import RandomForestR
8   # define dataset
9   X, y = make_regression(n_samples=1000, n_f                      m_s
10  # define the model
11  model = RandomForestRegressor()
12  # evaluate the model
13  cv = RepeatedKFold(n_splits=10, n_repeats=
14  n_scores = cross_val_score(model, X, y, sc               =-1
15  # report performance
16  print('MAE: %.3f (%.3f)' % (mean(n_scores)
```

## Start Machine Learning ✕

You can master applied Machine Learning **without math or fancy degrees**. Find out how in this *free* and *practical* course.

```
Email Address
```

☐ I consent to receive information about services and special offers by email. For more information, see the Privacy Policy.

**START MY EMAIL COURSE**

Running the example reports the mean and standa

**Note**: Your results may vary given the stochastic n differences in numerical precision. Consider runnir                                                    age

outcome

In this case, we can see the random forest ensemble with default hyperparameters achieves a MAE of about 90.

```
1  MAE: -90.149 (7.924)
```

We can also use the random forest model as a final model and make predictions for regression.

If the random forest ensemble is fit on all available data, then the *predict()* function can be called to make predictions on new data.

The example below demonstrates this on our regression dataset.

```
1  # random forest for making predictions for regression
2  from sklearn.datasets import make_regression
3  from sklearn.ensemble import RandomForestRegressor
4  # define dataset
5  X, y = make_regression(n_samples=1000, n_features=20, n_informative=15, noise=0.1, random_s
6  # define the model
7  model = RandomForestRegressor()
8  # fit the model on the whole dataset
9  model.fit(X, y)
10 # make a single prediction
11 row = [[-0.89483109,-1.0670149,-0.25448694,-0.53850126,0.21082105,1.37435592,0.71203659,0.7
12 yhat = model.predict(row)
13 print('Prediction: %d' % yhat[0])
```

Running the example fits the random forest ensemble model on the entire dataset and is then used to make a prediction on a new row of data, as we might when using the model in an application.

```
1  Prediction: -173
```

Now that we are familiar with using the scikit-learn API to evaluate and use random forest ensembles, let's look at configuring the model.

# Random Forest Hyperparameters

In this section, we will take a closer look at some of the hyperparameters you should consider tuning for the random forest ensemble and their effect on model performance.

## Explore Number of Samples

Each decision tree in the ensemble is fit on a boots

This can be turned off by setting the "*bootstrap*" ar                                                          hole training dataset will be used to train each decision

The "*max_samples*" argument can be set to a float                          size of the training dataset to make the bootstrap s

For example, if the training dataset has 100 rows,                          d each decision tree will be fit on a bootstrap sample

A smaller sample size will make trees more different, and a larger sample size will make the trees more similar. Setting *max_samples* to "*None*" will make the sample size the same size as the training dataset and this is the default.

The example below demonstrates the effect of different bootstrap sample sizes from 10 percent to 100 percent for the random forest algorithm.

```python
1   # explore random forest bootstrap sample size on performance
2   from numpy import mean
3   from numpy import std
4   from numpy import arange
5   from sklearn.datasets import make_classification
6   from sklearn.model_selection import cross_val_score
7   from sklearn.model_selection import RepeatedStratifiedKFold
8   from sklearn.ensemble import RandomForestClassifier
9   from matplotlib import pyplot
10
11  # get the dataset
12  def get_dataset():
13      X, y = make_classification(n_samples=1000, n_features=20, n_informative=15, n_redundant=
14      return X, y
15
16  # get a list of models to evaluate
17  def get_models():
18      models = dict()
19      # explore ratios from 10% to 100% in 10% increments
20      for i in arange(0.1, 1.1, 0.1):
21          key = '%.1f' % i
22          # set max_samples=None to use 100%
23          if i == 1.0:
24              i = None
25          models[key] = RandomForestClassifier(max_samples=i)
26      return models
27
28  # evaluate a given model using cross-validation
29  def evaluate_model(model, X, y):
30      # define the evaluation procedure
31      cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
32      # evaluate the model and collect the results
33      scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
34      return scores
35
36  # define dataset
37  X, y = get_dataset()
38  # get the models to evaluate
39  models = get_models()
40  # evaluate the models and store results
41  results, names = list(), list()
42  for name, model in models.items():
43      # evaluate the model
44      scores = evaluate_model(model, X, y)
45      # store the results
46      results.append(scores)
47      names.append(name)
48      # summarize the performance along the
49      print('>%s %.3f (%.3f)' % (name, mean(
50  # plot model performance for comparison
51  pyplot.boxplot(results, labels=names, show
52  pyplot.show()
```

Running the example first reports the mean accura

**Note**: Your results may vary given the stochastic n
differences in numerical precision. Consider runnir
outcome.

**Start Machine Learning**

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about services and special offers by email. For more information, see the Privacy Policy.

START MY EMAIL COURSE

In this case, the results suggest that using a bootstrap sample size that is equal to the size of the training dataset achieves the best results on this dataset.

This is the default and should probably be used in most cases.

```
1  >10 0.856 (0.031)
2  >20 0.873 (0.029)
3  >30 0.881 (0.021)
4  >40 0.891 (0.033)
5  >50 0.893 (0.025)
6  >60 0.897 (0.030)
7  >70 0.902 (0.024)
8  >80 0.903 (0.024)
9  >90 0.900 (0.026)
10 >100 0.903 (0.027)
```
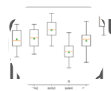
A box and whisker plot is created for the distribution of accuracy scores for each bootstrap sample size.

In this case, we can see a general trend that the larger the sample, the better the performance of the model.

You might like to extend this example and see what happens if the bootstrap sample size is larger or even much larger than the training dataset (e.g. you can set an integer value as the number of samples instead of a float percentage of the training dataset size).



Box Plot of Random Forest Bootstrap

## Explore Number of Features

The number of features that is randomly sampled for each split point is perhaps the most important feature to configure for random forest.

It is set via the *max_features* argument and defaults to the square root of the number of input features. In this case, for our test dataset, this would be *sqrt(20)* or about four features.

The example below explores the effect of the number of features randomly selected at each split point on model accuracy. We will try values from 1 to 7 and would expect a small value, around four, to perform well based on the heuristic.

```python
1  # explore random forest number of features effect on performance
2  from numpy import mean
3  from numpy import std
4  from sklearn.datasets import make_classification
5  from sklearn.model_selection import cross_val_score
6  from sklearn.model_selection import RepeatedStratifiedKFold
7  from sklearn.ensemble import RandomForestClassifier
8  from matplotlib import pyplot
9
10 # get the dataset
11 def get_dataset():
12     X, y = make_classification(n_samples=1000, n_features=20, n_informative=15, n_redundant=
13     return X, y
14
15 # get a list of models to evaluate
16 def get_models():
17     models = dict()
18     # explore number of features from 1 to 7
19     for i in range(1,8):
20         models[str(i)] = RandomForestClassifier(max_features=i)
21     return models
22
23 # evaluate a given model using cross-validation
24 def evaluate_model(model, X, y):
25     # define the evaluation procedure
26     cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
27     # evaluate the model and collect the results
28     scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
29     return scores
30
31 # define dataset
32 X, y = get_dataset()
33 # get the models to evaluate
34 models = get_models()
35 # evaluate the models and store results
36 results, names = list(), list()
37 for name, model in models.items():
38     # evaluate the model
39     scores = evaluate_model(model, X, y)
40     # store the results
41     results.append(scores)
42     names.append(name)
43     # summarize the performance along the
44     print('>%s %.3f (%.3f)' % (name, mean(
45 # plot model performance for comparison
46 pyplot.boxplot(results, labels=names, show
47 pyplot.show()
```

Running the example first reports the mean accura

**Note**: Your results may vary given the stochastic nature of the algorithm or evaluation procedure, or differences in numerical precision. Consider running the example a few times and compare the average outcome.

In this case, the results suggest that a value between three and five would be appropriate, confirming the sensible default of four on this dataset. A value of five might even be better given the smaller standard deviation in classification accuracy as compared to a value of three or four.

```
1  >1 0.897 (0.023)
2  >2 0.900 (0.028)
3  >3 0.903 (0.027)
4  >4 0.903 (0.022)
5  >5 0.903 (0.019)
6  >6 0.898 (0.025)
7  >7 0.900 (0.024)
```

A box and whisker plot is created for the distribution of accuracy scores for each feature set size.

We can see a trend in performance rising and peaking with values between three and five and falling again as larger feature set sizes are considered.



Box Plot of Random Forest Featu...

## Explore Number of Trees

The number of trees is another key hyperparamete

Typically, the number of trees is increased until the model performance stabilizes. Intuition might suggest that more trees will lead to overfitting, although this is not the case. Both bagging and random forest algorithms appear to be somewhat immune to overfitting the training dataset given the stochastic nature of the learning algorithm.

The number of trees can be set via the "*n_estimators*" argument and defaults to 100.

The example below explores the effect of the number of trees with values between 10 to 1,000.

```python
1   # explore random forest number of trees effect on performance
2   from numpy import mean
3   from numpy import std
4   from sklearn.datasets import make_classification
5   from sklearn.model_selection import cross_val_score
6   from sklearn.model_selection import RepeatedStratifiedKFold
7   from sklearn.ensemble import RandomForestClassifier
8   from matplotlib import pyplot
9   # get the dataset
10  def get_dataset():
11      X, y = make_classification(n_samples=1000, n_features=20, n_informative=15, n_redundant=
12      return X, y
13
14  # get a list of models to evaluate
15  def get_models():
16      models = dict()
17      # define number of trees to consider
18      n_trees = [10, 50, 100, 500, 1000]
19      for n in n_trees:
20          models[str(n)] = RandomForestClassifier(n_estimators=n)
21      return models
22
23  # evaluate a given model using cross-validation
24  def evaluate_model(model, X, y):
25      # define the evaluation procedure
26      cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
27      # evaluate the model and collect the results
28      scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
29      return scores
30
31  # define dataset
32  X, y = get_dataset()
33  # get the models to evaluate
34  models = get_models()
35  # evaluate the models and store results
36  results, names = list(), list()
37  for name, model in models.items():
38      # evaluate the model
39      scores = evaluate_model(model, X, y)
40      # store the results
41      results.append(scores)
42      names.append(name)
43      # summarize the performance along the
44      print('>%s %.3f (%.3f)' % (name, mean(
45  # plot model performance for comparison
46  pyplot.boxplot(results, labels=names, show
47  pyplot.show()
```

Running the example first reports the mean accura

**Note**: Your results may vary given the stochastic n
differences in numerical precision. Consider runnir
outcome.

**Start Machine Learning**   ✕

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about

services and special offers by email. For more

information, see the Privacy Policy.

START MY EMAIL COURSE

In this case, we can see that performance rises and stays flat after about 100 trees. Mean accuracy scores fluctuate across 100, 500, and 1,000 trees and this may be statistical noise.
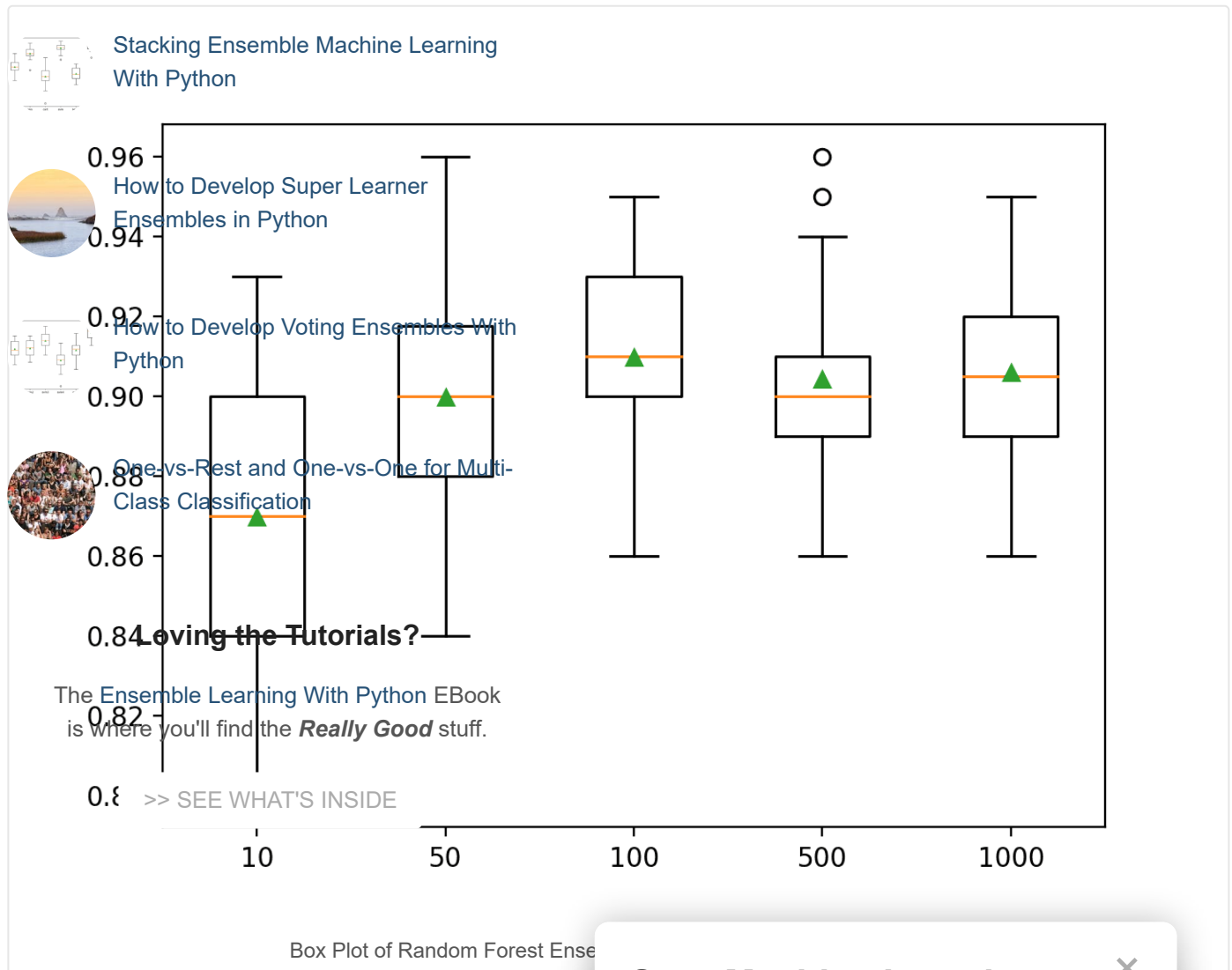
```
1  >10 0.870 (0.036)
2  >50 0.900 (0.028)
3  >100 0.910 (0.024)
4  >500 0.904 (0.024)
5  >1000 0.906 (0.023)
```

A box and whisker plot is created for the distribution of accuracy scores for each configured number of trees.



Box Plot of Random Forest Ense...

## Explore Tree Depth

A final interesting hyperparameter is the maximum

By default, trees are constructed to an arbitrary de although we can also explore fitting trees with diffe

The maximum tree depth can be specified via the maximum depth) by default.

The example below explores the effect of random

```
1  # explore random forest tree depth effect
```

```
 2  from numpy import mean
 3  from numpy import std
 4  from sklearn.datasets import make_classification
 5  from sklearn.model_selection import cross_val_score
 6  from sklearn.model_selection import RepeatedStratifiedKFold
 7  from sklearn.ensemble import RandomForestClassifier
 8  from matplotlib import pyplot
 9
10  # get the dataset
11  def get_dataset():
12      X, y = make_classification(n_samples=1000, n_features=20, n_informative=15, n_redundant=
13      return X, y
14
15  # get a list of models to evaluate
16  def get_models():
17      models = dict()
18      # consider tree depths from 1 to 7 and None=full
19      depths = [i for i in range(1,8)] + [None]
20      for n in depths:
21          models[str(n)] = RandomForestClassifier(max_depth=n)
22      return models
23
24  # evaluate a given model using cross-validation
25  def evaluate_model(model, X, y):
26      # define the evaluation procedure
27      cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)
28      # evaluate the model and collect the results
29      scores = cross_val_score(model, X, y, scoring='accuracy', cv=cv, n_jobs=-1)
30      return scores
31
32  # define dataset
33  X, y = get_dataset()
34  # get the models to evaluate
35  models = get_models()
36  # evaluate the models and store results
37  results, names = list(), list()
38  for name, model in models.items():
39      # evaluate the model
40      scores = evaluate_model(model, X, y)
41      # store the results
42      results.append(scores)
43      names.append(name)
44      # summarize the performance along the way
45      print('>%s %.3f (%.3f)' % (name, mean(scores), std(scores)))
46  # plot model performance for comparison
47  pyplot.boxplot(results, labels=names, showmeans=True)
48  pyplot.show()
```

Running the example first reports the mean accuracy for each configured maximum tree depth.

**Note**: Your results may vary given the stochastic n            differences in numerical precision. Consider runnin                                        age outcome.

In this case, we can see that larger depth results in maximum depth achieving the best performance o

```
1  >1 0.771 (0.040)
2  >2 0.807 (0.037)
3  >3 0.834 (0.034)
4  >4 0.857 (0.030)
5  >5 0.872 (0.025)
6  >6 0.887 (0.024)
7  >7 0.890 (0.025)
8  >None 0.903 (0.027)
```

**Start Machine Learning**    ×

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about
services and special offers by email. For more
information, see the Privacy Policy.

START MY EMAIL COURSE

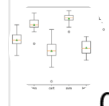A box and whisker plot is created for the distribution of accuracy scores for each configured maximum tree depth.

In this case, we can see a trend of improved performance with increase in tree depth, supporting the default of no maximum depth.
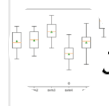
Random Forest Maximum Tree Depth vs. Classification Accuracy

# Common Questions

In this section we will take a closer look at some com om forest ensemble procedure.

**Q. What algorithm should be used in the ensem**

Random forest is designed to be an ensemble of d

**Q. How many ensemble members should be us**

The number of trees should be increased until no f ur dataset.

> *As a starting point, we suggest using at lea*
> *profiles are still improving at 1,000 trees, th*

* Random Forests, 2001.

**Never miss a tutorial:**

**Books**



* Applied Predictive Modeling, 2013.
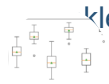* The Elements of Statistical Learning, 2016.

**Picked for you:**

* An Introduction to Statistical Learning with Applications in R, 2014.

How to Develop Multi-Output Regression Models with Python

* sklearn.ensemble.RandomForestRegressor API.
* sklearn.ensemble.RandomForestClassifier API.
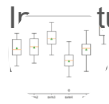
Stacking Ensemble Machine Learning With Python

## Articles

* Random Forest, Wikipedia.

How to Develop Super Learner Ensembles in Python

## Summary

In this tutorial, you discovered how to develop random forest ensembles for classification and regression.

How to Develop Voting Ensembles With Python

Specifically, you learned:

One-vs-Rest and One-vs-One for Multi-Class Classification ensemble is an ensemble of decision trees and a natural extension of bagging.

* How to use the random forest ensemble for classification and regression with scikit-learn.
* How to explore the effect of random forest model hyperparameters on model performance.

**Do you have any questions?**
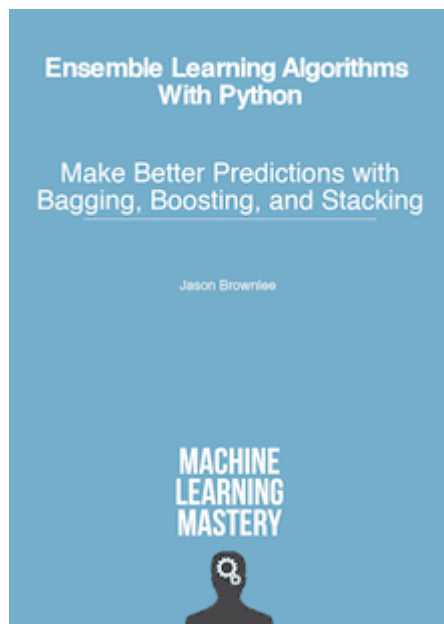
**Loving the Tutorials?**

Ask your questions in the comments below and I will do my best to answer.

The Ensemble Learning With Python EBook is where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE

# Get a Handle on Modern Ensemble Learning!

**Ensemble Learning Algorithms With Python**

Make Better Predictions with Bagging, Boosting, and Stacking

Jason Brownlee

**MACHINE LEARNING MASTERY**

**Improve Your Predictions in Minutes**

It provid...

*Stacking*, *Voting*, E... ...ore...

**Bring M...**

**Start Machine Learning** ✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about services and special offers by email. For more information, see the Privacy Policy.

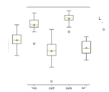**START MY EMAIL COURSE**

**Never miss a tutorial:**
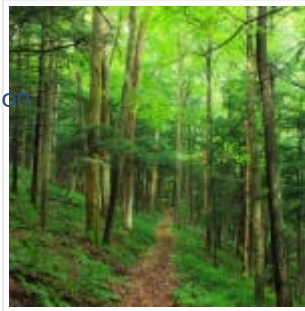
# More On This Topic

**Picked for you:**

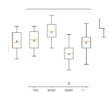How to Develop Multi-Output Regression
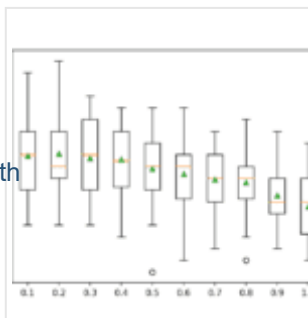Models with Python

Stacking Ensemble Machine Learning
With Python

Bagging and Random Forest Ensemble Algorithms for…

How to Develop Super Learner
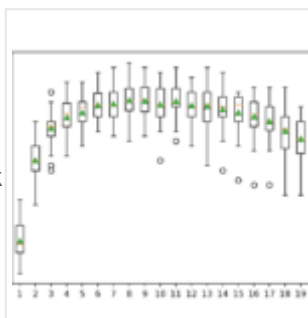Ensembles in Python

How to Develop Voting Ensembles With
Python

How to Develop Random Forest Ensembles With XGBoost

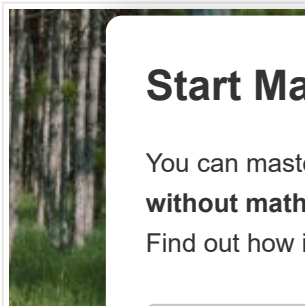One-vs-Rest and One-vs-One for Multi-
Class Classification

### Loving the Tutorials?

The Ensemble Learning With Python EBook
is where you'll find the **Really Good** stuff.

>> SEE WHAT'S INSIDE

How to Develop a Random Subspace Ensemble With Python

## Start Machine Learning

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about
services and special offers by email. For more
information, see the Privacy Policy.

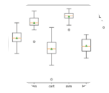**START MY EMAIL COURSE**

How to Implement Rando

**Never miss a tutorial:**



**Picked for you:**

How to Develop Multi-Output Regression Models with Python
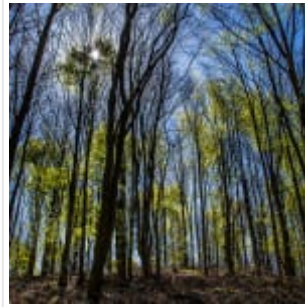
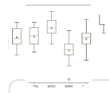Use Random Forest: Testing 179 Classifiers on 121 Datasets

Stacking Ensemble Machine Learning With Python

How to Develop Super Learner Ensembles in Python

Tune Machine Learning Algorithms in R (random forest…

How to Develop Voting Ensembles With Python

**About Jason Brownlee**

One-vs-Rest and One-vs-One for Multi-Class Classification

Jason Brownlee, PhD is a machine learning specialist who teaches developers how to get results with modern machine learning methods via hands-on tutorials.

View all posts by Jason Brownlee →

**Loving the Tutorials?**

‹ How to Develop Voting Ensembles With Python

How to Develop an Extra Trees Ensemble with Python ›

The Ensemble Learning With Python EBook is where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE

## 36 Responses to *How to Develop a Random Forest Ensemble in Python*

**dcart** April 20, 2020 at 6:20 am #

Hi Jason

Hope you are doing well in this time of lock down.

Great articles as usual. Although I like to apply RF memory as the data is huge.

Any advise using RF for huge data set?

Thanks
Dennis

**Start Machine Learning** ✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about services and special offers by email. For more information, see the Privacy Policy.

START MY EMAIL COURSE

**Never miss a tutorial:**

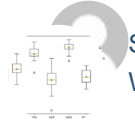**Jason Brownlee** April 20, 2020 at 7:36 am # REPLY

Perhaps prepare a prototype on a small sample of data first to see if it is effective.

More ideas here:

**Picked for you:** https://machinelearningmastery.com/faq/single-faq/how-can-i-run-large-models-or-models-on-lots-of-data

How to Develop Multi-Output Regression Models with Python

**Elie** April 23, 2020 at 4:05 am # REPLY
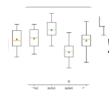
Stacking Ensemble Machine Learning With Python Hi Jason, Are you planning a new book on EnsembleS?

How to Develop Super Learner Ensembles in Python

**Jason Brownlee** April 23, 2020 at 6:12 am # REPLY

Not sure at this stage.

How to Develop Voting Ensembles With Python

Are there ensemble topics you'd like me to write about?

One-vs-Rest and One-vs-One for Multi-Class Classification

**Elie** April 24, 2020 at 5:17 am # REPLY

probably advanced stacking and how to win kaggle/data science competitions

**Loving the Tutorials?**

The Ensemble Learning With Python EBook **Jason Brownlee** April 24, 2020 at 5:54 am # REPLY
is where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE

https://machinelearningmastery.com/stacking-ensemble-machine-learning-with-python/

And this:

https://machinelearningmastery.com/super-learner-ensemble-in-python/

**Stelios** April 26, 2020 at 4:34 pm #

Thank you so much for your great support

**Start Machine Learning** ✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about services and special offers by email. For more information, see the Privacy Policy.

START MY EMAIL COURSE

**Jason Brownlee** April 27, 2020 at 5:30 a

You're welcome.

**Never miss a tutorial:**

**Picked for you:**

How to Develop Multi-Output Regression Models with Python

Stacking Ensemble Machine Learning With Python

How to Develop Super Learner Ensembles in Python

How to Develop Voting Ensembles With Python

One-vs-Rest and One-vs-One for Multi-Class Classification

**Loving the Tutorials?**

The Ensemble Learning With Python EBook is where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE

REPLY ↩

**David Sanchez** May 10, 2020 at 11:25 pm #
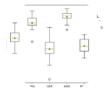
I'm implementing a Random Forest and I'm getting a shifted time-series in the predictions. If I build the model for predicting e.g. 4 steps ahead, my time-series of predictions seems 4 steps shifted to the right comparing to my time-series of observations. If I try to predict 16 steps ahead, it seems 16 steps shifted.

Any idea why this could be happening?

Thanks for all your tutorials!

REPLY ↩

**Jason Brownlee** May 11, 2020 at 6:00 am #

Yes, it sounds like the model has learned a persistence (no skill) forecast. E.g. it predicts the input as the output.

REPLY ↩

**David Sanchez** May 14, 2020 at 6:28 pm #

Hi Jason,

Thanks for your reply.

I'm also wondering, if I try to build a model where my train set has more variables than my test set, how should I proceed?

As far as I've seen about it, I should recreate those missing variables in my test dataframe and set them as 0.

REPLY ↩

**Jason Brownlee** May 15, 2020 at 5:57 am #

The number of variables (columns) must be the same in train and test sets.

**Grzegorz Kępisty** May 26, 2020 at 10:46 pm

Very nice tutorial of RF usage!
It is really practical to know good practices on thos very competitive in real industrial applications! (ofte Networks).
Regards!

**Jason Brownlee** May 27, 2020 at 7:54 a

**Start Machine Learning** ✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about services and special offers by email. For more information, see the Privacy Policy.

START MY EMAIL COURSE

Thanks!

**Never miss a tutorial:**

Agreed. XGBoost more so.



**Picked for you:**

**Dilo** May 28, 2020 at 10:03 am #

Hi Jason,

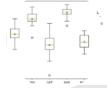[How to Develop Multi-Output Regression Models with Python]

ase, check It.

"This means that larger negative MAE are better and a perfect model has a MAE of 0."

Thank you! enjoying a lot this stuff.

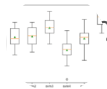[Stacking Ensemble Machine Learning With Python]

**Jason Brownlee** May 28, 2020 at 1:24 pm #

[How to Develop Super Learner Ensembles in Python]

It is correct.

-10 is greater than -100.

[How to Develop Voting Ensembles With Python]

0 is greater than -10.

[One-vs-Rest and One-vs-One for Multi-Class Classification]

**manuela** October 16, 2020 at 4:26 pm #

Hello Jason, Please I have a question

I have the following situation that is already programmed with Logistic regression, I have tried the same program with Random Forest in order to check how it could improve the accuracy.

Actually, the accuracy was improved, but I don't know if it is logical to use the Random Forest in my problem case.

My case study is as follow :

Based                                           to predict if a customer will buy a product or not depending on his prior history. i.e to know how much a customer bought the same product previously, and how much he just check it without buying it

The used data has the following structure:

Id clients CurrectProd P1+ P1- P2+ P2- P3+ P3- ..

10 CL1 P1, P3 6 1 0 0 8 2 0 0 1

11 CL1 P1, P2 7 1 5 2 0 0 0 0 1

with:

CurrentProd: means a list of products that I need t

P1+: mean how many time à client buy product 1,

P1-: refers to the number that a client checked a pr

columns present all products existing in the market

PRODUCT) and at each row the most of those row

CurrentPRod

So I want to know if the random forest could be use

PS: I must use the data as it is without any change in features or structure

**Never miss a tutorial:**

**Picked for you:**

**manuela** October 16, 2020 at 4:43 pm #

REPLY

How to Develop Multi-Output Regression Models with Python

```
Id..|..clients..|..CurrectProd..|.P1+.|.P1-.|.P2+.|.P2-.|.P3+.|.P3-.| … .|.PN+.|.PN-.|.Output
10|…CL1…|…P1,P3…|…6…|…1…|…0….|…0…|..8…|…2…| … .|…0…|…0…|….1
11|…CL1…|…P1, P2……|…7….|..1…|…5….|…2…|…0…|..0…|. … .|…0…|…0…|….1
```

Stacking Ensemble Machine Learning with Python

**Jason Brownlee** October 17, 2020 at 5:58 am #

REPLY

Perhaps try it and compare results.

The key will to find an appropriate representation for the problem. This may give you ideas (replace site with product):

https://machinelearningmastery.com/faq/single-faq/how-to-develop-forecast-models-for-multiple-sites

How to Develop Super Learner Ensembles in Python

How to Develop Voting Ensembles With Python

**manuela** October 17, 2020 at 6:45 am #

REPLY

Thanks for your quick replay

I have already tried it, and it gives me a good result,

but I want to know if it is logical to use it with 200 features (Product1, Product2….)

One-vs-Rest and One-vs-One for Multi-Class Classification

**Loving the Tutorials?**

The Ensemble Learning With Python EBook is where you find the *Really Good* stuff.

**Jason Brownlee** October 17, 2020 at 1:41 pm #

REPLY

Use the features that result in the best performance, regardless of how many.

>> SEE WHAT'S INSIDE

**dinesh** December 21, 2020 at 7:29 pm #

REPLY

how to decide these paramters

n_samples=1000, n_features=20, n_informative=1

any suggestion on this please.

**Jason Brownlee** December 22, 2020 at 6:

This defines the test problem, it is con

The ideas is you replace this with your own dat

## Never miss a tutorial:

**Picked for you:**

How to Develop Multi-Output Regression Models with Python

Stacking Ensemble Machine Learning With Python

How to Develop Super Learner Ensembles in Python

How to Develop Voting Ensembles With Python

One-vs-Rest and One-vs-One for Multi-Class Classification

**Loving the Tutorials?**

The Ensemble Learning With Python EBook is where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE

---

**dinesh** December 25, 2020 at 1:51 pm #　　　　　　REPLY ↩

what are the best practice , also if i want to learn about the meaning of these parameter. please share some information about the hyper parameter tuning.

**Jason Brownlee** December 26, 2020 at 5:08 am #　　　REPLY ↩

Best practice for test problems? I don't understand.

Best practice for hyperparameter tuning is in the above example, e.g. grid search.

**Jim** February 17, 2021 at 6:03 pm #　　　　　　REPLY ↩

Perfect!

**Jason Brownlee** February 18, 2021 at 5:12 am #　　　REPLY ↩

Thanks!

**Alex** September 8, 2021 at 12:29 am #　　　　　　REPLY ↩

Thanks for such comprehensive article,

I wonder to know is there any way to find out that under which condition my model has wrong prediction, is there any way to find (range) values for features that tell me that the model is not reliable. Or let machine learn that when the prediction could not be reliable.

I think it could be, some how , the other way around of machine learning ,isnt it?

Any suggestions for it?

**Adrian Tam** September 8, 2021 at [...]

Quite impossible to know because [...] data provided. You're simply asking wh[...]

**Francisco Pérez Liébana** March 12, 2021 [...]

Hello Jason,

---

**Start Machine Learning**　　✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about services and special offers by email. For more information, see the Privacy Policy.

START MY EMAIL COURSE

Thanks for your articles, they are very useful!

**Never miss a tutorial:**

Do you know how can I get a graphic representation of the trees in the trained model ? I was trying to use export_graphviz sklearn but using "cross_val_scores" function fitting estimator on its own, i don´t know how to use export_gaphviz function.

Thanks in advance for your answer.

**Picked for you:**

Francisco

How to Develop Multi-Output Regression Models with Python

**Jason Brownlee** March 13, 2021 at 5:33 am # REPLY

Stacking Ensemble Machine Learning With Python believe it's possible but I have not done it before, sorry Francisco.

How to Develop Super Learner Ensembles in Python **Ilan Sharfer** May 12, 2021 at 11:59 pm # REPLY

Hi Jason,

How to Develop Voting Ensembles With Python
Thanks for the clear and useful introduction.

I have a question on how the Random Forest algorithm handles missing features.

For example suppose the data set is a 24H time series, for which I want to build a classifier.
One-vs-Rest and One-vs-One for Multi-Class Classification
Some of the features are available only in daytime, some only in night-time, and some others are partly unavailable.

What is the best way to adapt the algorithm to address this task.

**Loving the Tutorials?**
Thanks in advance,

The Ensemble Learning With Python EBook
Ilan is where you'll find the *Really Good* stuff.

—— >> SEE WHAT'S INSIDE ——

**Jason Brownlee** May 13, 2021 at 6:03 am # REPLY

Perhaps try a suite of approaches for handling the missing data and discover what works well or best for your dataset.

**Ahmad Afif Aulia Hariz** July 13, 2021 at 10

Thanks a lot for the article.

I've just build my own RF Regressor, i have (2437, about 0.7

I want to improve it into 0.95. Any suggestion?
Thanks for help!

**Never miss a tutorial:**

**Jason Brownlee** July 14, 2021 at 5:29 am #

Perhaps some of these ideas will help:
https://machinelearningmastery.com/machine-learning-performance-improvement-cheat-sheet/

**Picked for you:**

How to Develop Multiple Regression
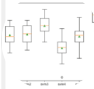Models with Python

**ANANTHAKRISHNAN** July 14, 2021 at 7:41 pm #

Hai sir,

I'm having two seperate data frames. One for Training and another one for Testing. I need to perform Random Forest Classification. How to load these files to Random Forest without splitting.

Stacking Ensemble Machine Learning
With Python

How to Develop Super Learner
Ensembles in Python

**Leave a Reply**

How to Develop Voting Ensembles With
Python

One-vs-Rest and One-vs-One for Multi-
Class Classification

**Loving the Tutorials?**

Name (required)

The Ensemble Learning With Python EBook
is where you'll find the *Really Good* stuff.

Email (will not be published) (required)

>> SEE WHAT'S INSIDE

Website

SUBMIT COMMENT

**Welcome!**
I'm *Jason Brownlee* PhD
and I **help developers** get results with
Read more

**Start Machine Learning** ✕

You can master applied Machine Learning
**without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about

services and special offers by email. For more

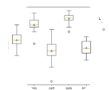information, see the Privacy Policy.

START MY EMAIL COURSE

**Never miss a tutorial:**
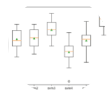
**Picked for you:**

How to Develop Multi-Output Regression Models with Python

Stacking Ensemble Machine Learning With Python

How to Develop Super Learner Ensembles in Python

How to Develop Voting Ensembles With Python

One-vs-Rest and One-vs-One for Multi-Class Classification

**Loving the Tutorials?**

The Ensemble Learning With Python EBook is where you'll find the *Really Good* stuff.

>> SEE WHAT'S INSIDE

© 2021 Machine Learning Mastery. All Rights Reserved.
LinkedIn | Twitter | Facebook | Newsletter | RSS

Privacy | Disclaimer | Terms | Contact | Sitemap | Search

**Start Machine Learning** ✕

You can master applied Machine Learning **without math or fancy degrees**.
Find out how in this *free* and *practical* course.

Email Address

☐ I consent to receive information about services and special offers by email. For more information, see the Privacy Policy.

START MY EMAIL COURSE