

[Start Here](#)[Blog](#)[Books](#)[About](#)[Contact](#)

Need help with machine learning? [Take the FREE Crash-Course.](#)

# How to Implement Bagging From Scratch With Python

by **Jason Brownlee** on November 11, 2016 in **Algorithms From Scratch**



Decision trees are a simple and powerful predictive modeling technique, but they suffer from high-variance.

This means that trees can get very different results given different training data.

A technique to make decision trees more robust and to achieve better performance is called bootstrap aggregation or bagging for short.

In this tutorial, you will discover how to implement the bagging procedure with decision trees from scratch with Python.

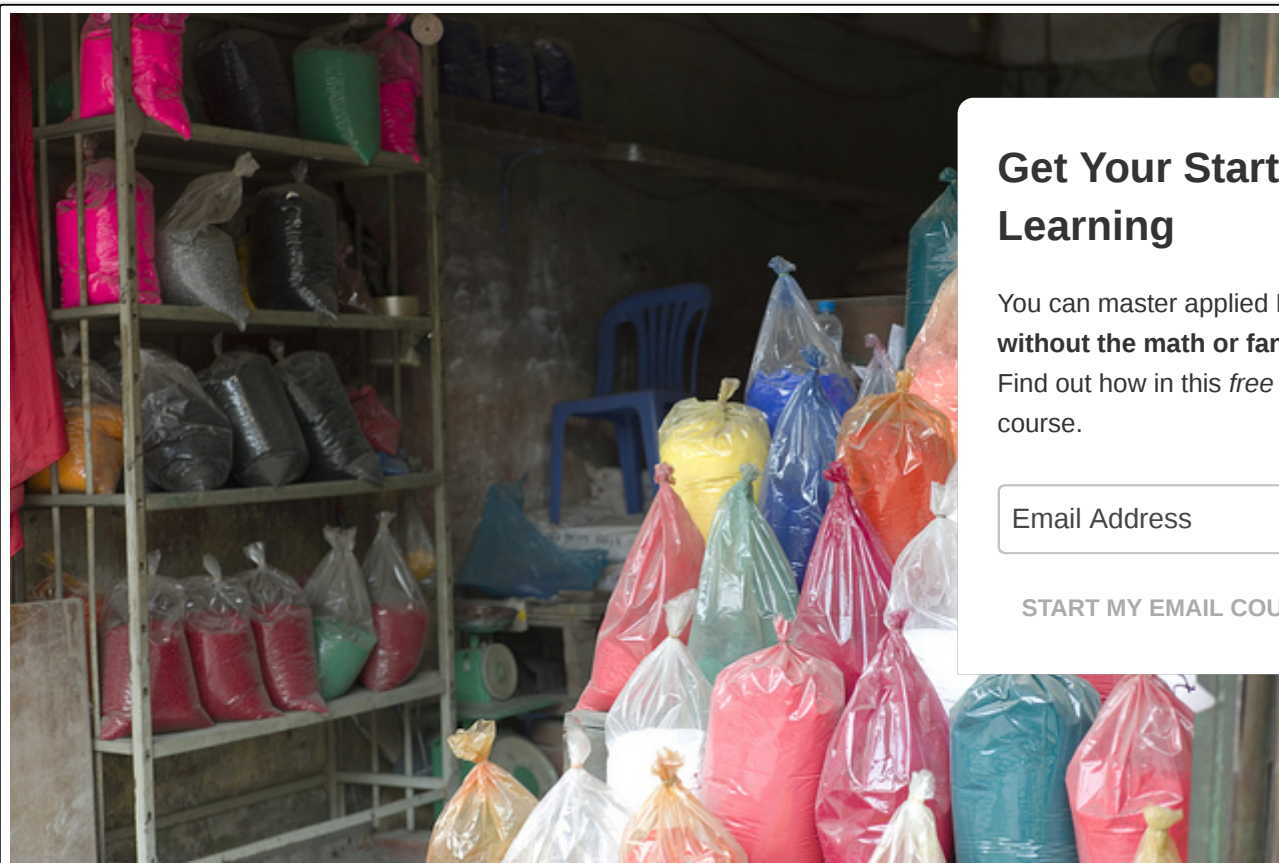
After completing this tutorial, you will know:

[Get Your Start in Machine Learning](#)

- How to create a bootstrap sample of your dataset.
- How to make predictions with bootstrapped models.
- How to apply bagging to your own predictive modeling problems.

Let's get started.

- **Update Jan/2017:** Changed the calculation of `fold_size` in `cross_validation_split()` to always be an integer. Fixes issues with Python 3.
- **Update Feb/2017:** Fixed a bug in `build_tree`.
- **Update Aug/2017:** Fixed a bug in Gini calculation, added the missing weighting of group Gini scores by group size (thanks Michael!).



How to Implement Bagging From Scratch With Python  
Photo by [Michael Cory](#), some rights reserved.

## Get Your Start in Machine Learning

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

START MY EMAIL COURSE

Get Your Start in Machine Learning

## Descriptions

This section provides a brief description to Bootstrap Aggregation and the Sonar dataset that will be used in this tutorial.

### Bootstrap Aggregation Algorithm

A bootstrap is a sample of a dataset with replacement.

This means that a new dataset is created from a random sample of an existing dataset where a given row may be selected and added more than once to the sample.

It is a useful approach to use when estimating values such as the mean for a broader dataset, when you only have a limited dataset available. By creating samples of your dataset and estimating the mean from those samples, you can take the average of the means to get a better estimate of the true mean of the underlying problem.

This same approach can be used with machine learning algorithms that have a high variance, such as decision trees. By creating multiple bootstrap samples of data and the average output of those models used to make predictions. This is known as bagging for short.

Variance means that an algorithm's performance is sensitive to the training data, with high variance. When the training data is changed, the more the performance of the algorithm will vary.

The performance of high variance machine learning algorithms like unpruned decision trees can be improved by averaging their predictions. Results are often better than a single decision tree.

Another benefit of bagging in addition to improved performance is that the bagged decision trees can be added until a maximum in performance is achieved.

### Sonar Dataset

The dataset we will use in this tutorial is the Sonar dataset.

This is a dataset that describes sonar chirp returns bouncing off different surfaces. The 60 input variables are the strength of the returns at different angles. It is a binary classification problem that requires a model to differentiate rocks from metal cylinders.

### Get Your Start in Machine Learning

You can master applied Machine Learning **without the math or fancy degree.**

Find out how in this *free* and *practical* email course.

START MY EMAIL COURSE

It is a well-understood dataset. All of the variables are continuous and generally in the range of 0 to 1. The output variable is a string "M" for mine and "R" for rock, which will need to be converted to integers 1 and 0.

By predicting the class with the most observations in the dataset (M or mines) the Zero Rule Algorithm can achieve an accuracy of 53%.

You can learn more about this dataset at the [UCI Machine Learning repository](#).

Download the dataset for free and place it in your working directory with the filename **sonar.all-data.csv**.

## Tutorial

This tutorial is broken down into 2 parts:

1. Bootstrap Resample.
2. Sonar Dataset Case Study.

These steps provide the foundation that you need to implement and apply bootstrap aggregation with problems.

### 1. Bootstrap Resample

Let's start off by getting a strong idea of how the bootstrap method works.

We can create a new sample of a dataset by randomly selecting rows from the dataset and adding them until the number of rows or until the size of the new dataset matches a ratio of the size of the original dataset.

We can allow sampling with replacement by not removing the row that was selected so that it is available for future selections.

Below is a function named **subsample()** that implements this procedure. The **randrange()** function from the random module is used to select a random row index to add to the sample each iteration of the loop. The default size of the sample is the size of the original dataset.

```
1 # Create a random subsample from the dataset with replacement
2 def subsample(dataset, ratio=1.0):
3     sample = list()
4     n_sample = round(len(dataset) * ratio)
5     while len(sample) < n_sample:
```

### Get Your Start in Machine Learning

You can master applied Machine Learning **without the math or fancy degree.**

Find out how in this *free* and *practical* email course.

START MY EMAIL COURSE

Get Your Start in Machine Learning

```

6         index = randrange(len(dataset))
7         sample.append(dataset[index])
8     return sample

```

We can use this function to estimate the mean of a contrived dataset.

First, we can create a dataset with 20 rows and a single column of random numbers between 0 and 9 and calculate the mean value.

We can then make bootstrap samples of the original dataset, calculate the mean, and repeat this process until we have a list of means. Taking the average of these sample means should give us a robust estimate of the mean of the entire dataset.

The complete example is listed below.

Each bootstrap sample is created as a 10% sample of the original 20 observation dataset (or 2 observations). We then experiment by creating 1, 10, 100 bootstrap samples of the original dataset, calculate their mean value, then average all of those estimates.

```

1  from random import seed
2  from random import random
3  from random import randrange
4
5
6  # Create a random subsample from the dataset with replacement
7  def subsample(dataset, ratio=1.0):
8      sample = list()
9      n_sample = round(len(dataset) * ratio)
10     while len(sample) < n_sample:
11         index = randrange(len(dataset))
12         sample.append(dataset[index])
13     return sample
14
15
16 # Calculate the mean of a list of numbers
17 def mean(numbers):
18     return sum(numbers) / float(len(numbers))
19
20
21 seed(1)
22 # True mean
23 dataset = [[randrange(10)] for i in range(20)]
24 print('True Mean: %.3f' % mean([row[0] for row in dataset]))
25 # Estimated means
26 ratio = 0.10
27 for size in [1, 10, 100]:

```

## Get Your Start in Machine Learning

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

START MY EMAIL COURSE

Get Your Start in Machine Learning

```

28 sample_means = list()
29 for i in range(size):
30     sample = subsample(dataset, ratio)
31     sample_mean = mean([row[0] for row in sample])
32     sample_means.append(sample_mean)
33 print('Samples=%d, Estimated Mean: %.3f' % (size, mean(sample_means)))

```

Running the example prints the original mean value we aim to estimate.

We can then see the estimated mean from the various different numbers of bootstrap samples. We can see that with 100 samples we achieve a good estimate of the mean.

```

1 True Mean: 4.450
2 Samples=1, Estimated Mean: 4.500
3 Samples=10, Estimated Mean: 3.300
4 Samples=100, Estimated Mean: 4.480

```

Instead of calculating the mean value, we can create a model from each subsample.

Next, let's see how we can combine the predictions from multiple bootstrap models.

## 2. Sonar Dataset Case Study

In this section, we will apply the Random Forest algorithm to the Sonar dataset.

The example assumes that a CSV copy of the dataset is in the current working directory with the file

The dataset is first loaded, the string values converted to numeric and the output column is converted to numeric. This is achieved with helper functions **load\_csv()**, **str\_column\_to\_float()** and **str\_column\_to\_int()** to load

We will use k-fold cross validation to estimate the performance of the learned model on unseen data. This means that we will construct and evaluate k models and estimate the performance as the mean model error. Classification accuracy will be used to evaluate each model. These behaviors are provided in the **cross\_validation\_split()**, **accuracy\_metric()** and **evaluate\_algorithm()** helper functions.

We will also use an implementation of the Classification and Regression Trees (CART) algorithm adapted for bagging including the helper functions **test\_split()** to split a dataset into groups, **gini\_index()** to evaluate a split point, **get\_split()** to find an optimal split point, **to\_terminal()**, **split()** and **build\_tree()** used to create a single decision tree, **predict()** to make a prediction with a decision tree and **bootstrap()** used in the previous step to make a subsample of the training dataset

### Get Your Start in Machine Learning

You can master applied Machine Learning **without the math or fancy degree.**

Find out how in this *free* and *practical* email course.

START MY EMAIL COURSE

Get Your Start in Machine Learning

A new function named **bagging\_predict()** is developed that is responsible for making a prediction with each decision tree and combining the predictions into a single return value. This is achieved by selecting the most common prediction from the list of predictions made by the bagged trees.

Finally, a new function named **bagging()** is developed that is responsible for creating the samples of the training dataset, training a decision tree on each, then making predictions on the test dataset using the list of bagged trees.

The complete example is listed below.

```
1 # Bagging Algorithm on the Sonar dataset
2 from random import seed
3 from random import randrange
4 from csv import reader
5
6 # Load a CSV file
7 def load_csv(filename):
8     dataset = list()
9     with open(filename, 'r') as file:
10         csv_reader = reader(file)
11         for row in csv_reader:
12             if not row:
13                 continue
14             dataset.append(row)
15     return dataset
16
17 # Convert string column to float
18 def str_column_to_float(dataset, column):
19     for row in dataset:
20         row[column] = float(row[column].strip())
21
22 # Convert string column to integer
23 def str_column_to_int(dataset, column):
24     class_values = [row[column] for row in dataset]
25     unique = set(class_values)
26     lookup = dict()
27     for i, value in enumerate(unique):
28         lookup[value] = i
29     for row in dataset:
30         row[column] = lookup[row[column]]
31     return lookup
32
33 # Split a dataset into k folds
34 def cross_validation_split(dataset, n_folds):
35     dataset_split = list()
36     dataset_copy = list(dataset)
```

## Get Your Start in Machine Learning ×

You can master applied Machine Learning **without the math or fancy degree.**

Find out how in this *free* and *practical* email course.

START MY EMAIL COURSE

Get Your Start in Machine Learning

```

37     fold_size = int(len(dataset) / n_folds)
38     for i in range(n_folds):
39         fold = list()
40         while len(fold) < fold_size:
41             index = randrange(len(dataset_copy))
42             fold.append(dataset_copy.pop(index))
43         dataset_split.append(fold)
44     return dataset_split
45
46 # Calculate accuracy percentage
47 def accuracy_metric(actual, predicted):
48     correct = 0
49     for i in range(len(actual)):
50         if actual[i] == predicted[i]:
51             correct += 1
52     return correct / float(len(actual)) * 100.0
53
54 # Evaluate an algorithm using a cross validation split
55 def evaluate_algorithm(dataset, algorithm, n_folds, *args):
56     folds = cross_validation_split(dataset, n_folds)
57     scores = list()
58     for fold in folds:
59         train_set = list(folds)
60         train_set.remove(fold)
61         train_set = sum(train_set, [])
62         test_set = list()
63         for row in fold:
64             row_copy = list(row)
65             test_set.append(row_copy)
66             row_copy[-1] = None
67         predicted = algorithm(train_set, test_set, *args)
68         actual = [row[-1] for row in fold]
69         accuracy = accuracy_metric(actual, predicted)
70         scores.append(accuracy)
71     return scores
72
73 # Split a dataset based on an attribute and an attribute value
74 def test_split(index, value, dataset):
75     left, right = list(), list()
76     for row in dataset:
77         if row[index] < value:
78             left.append(row)
79         else:
80             right.append(row)
81     return left, right
82
83 # Calculate the Gini index for a split dataset

```

## Get Your Start in Machine Learning

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

START MY EMAIL COURSE

Get Your Start in Machine Learning



```

84 def gini_index(groups, classes):
85     # count all samples at split point
86     n_instances = float(sum([len(group) for group in groups]))
87     # sum weighted Gini index for each group
88     gini = 0.0
89     for group in groups:
90         size = float(len(group))
91         # avoid divide by zero
92         if size == 0:
93             continue
94         score = 0.0
95         # score the group based on the score for each class
96         for class_val in classes:
97             p = [row[-1] for row in group].count(class_val) / size
98             score += p * p
99         # weight the group score by its relative size
100        gini += (1.0 - score) * (size / n_instances)
101    return gini
102
103 # Select the best split point for a dataset
104 def get_split(dataset):
105     class_values = list(set(row[-1] for row in dataset))
106     b_index, b_value, b_score, b_groups = 999, 999, 999, None
107     for index in range(len(dataset[0])-1):
108         for row in dataset:
109             # for i in range(len(dataset)):
110             #     row = dataset[randrange(len(dataset))]
111             groups = test_split(index, row[index], dataset)
112             gini = gini_index(groups, class_values)
113             if gini < b_score:
114                 b_index, b_value, b_score, b_groups = index, row[index], gini, groups
115     return {'index':b_index, 'value':b_value, 'groups':b_groups}
116
117 # Create a terminal node value
118 def to_terminal(group):
119     outcomes = [row[-1] for row in group]
120     return max(set(outcomes), key=outcomes.count)
121
122 # Create child splits for a node or make terminal
123 def split(node, max_depth, min_size, depth):
124     left, right = node['groups']
125     del(node['groups'])
126     # check for a no split
127     if not left or not right:
128         node['left'] = node['right'] = to_terminal(left + right)
129         return
130     # check for max depth

```

## Get Your Start in Machine Learning ×

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

START MY EMAIL COURSE

Get Your Start in Machine Learning

```

131     if depth >= max_depth:
132         node['left'], node['right'] = to_terminal(left), to_terminal(right)
133         return
134     # process left child
135     if len(left) <= min_size:
136         node['left'] = to_terminal(left)
137     else:
138         node['left'] = get_split(left)
139         split(node['left'], max_depth, min_size, depth+1)
140     # process right child
141     if len(right) <= min_size:
142         node['right'] = to_terminal(right)
143     else:
144         node['right'] = get_split(right)
145         split(node['right'], max_depth, min_size, depth+1)
146
147 # Build a decision tree
148 def build_tree(train, max_depth, min_size):
149     root = get_split(train)
150     split(root, max_depth, min_size, 1)
151     return root
152
153 # Make a prediction with a decision tree
154 def predict(node, row):
155     if row[node['index']] < node['value']:
156         if isinstance(node['left'], dict):
157             return predict(node['left'], row)
158         else:
159             return node['left']
160     else:
161         if isinstance(node['right'], dict):
162             return predict(node['right'], row)
163         else:
164             return node['right']
165
166 # Create a random subsample from the dataset with replacement
167 def subsample(dataset, ratio):
168     sample = list()
169     n_sample = round(len(dataset) * ratio)
170     while len(sample) < n_sample:
171         index = randrange(len(dataset))
172         sample.append(dataset[index])
173     return sample
174
175 # Make a prediction with a list of bagged trees
176 def bagging_predict(trees, row):
177     predictions = [predict(tree, row) for tree in trees]

```

## Get Your Start in Machine Learning ×

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

START MY EMAIL COURSE

Get Your Start in Machine Learning

```

178     return max(set(predictions), key=predictions.count)
179
180 # Bootstrap Aggregation Algorithm
181 def bagging(train, test, max_depth, min_size, sample_size, n_trees):
182     trees = list()
183     for i in range(n_trees):
184         sample = subsample(train, sample_size)
185         tree = build_tree(sample, max_depth, min_size)
186         trees.append(tree)
187     predictions = [bagging_predict(trees, row) for row in test]
188     return(predictions)
189
190 # Test bagging on the sonar dataset
191 seed(1)
192 # load and prepare data
193 filename = 'sonar.all-data.csv'
194 dataset = load_csv(filename)
195 # convert string attributes to integers
196 for i in range(len(dataset[0])-1):
197     str_column_to_float(dataset, i)
198 # convert class column to integers
199 str_column_to_int(dataset, len(dataset[0])-1)
200 # evaluate algorithm
201 n_folds = 5
202 max_depth = 6
203 min_size = 2
204 sample_size = 0.50
205 for n_trees in [1, 5, 10, 50]:
206     scores = evaluate_algorithm(dataset, bagging, n_folds, max_depth, min_size, sample_size)
207     print('Trees: %d' % n_trees)
208     print('Scores: %s' % scores)
209     print('Mean Accuracy: %.3f%%' % (sum(scores)/float(len(scores))))

```

A k value of 5 was used for cross-validation, giving each fold  $208/5 = 41.6$  or just over 40 records to

Deep trees were constructed with a max depth of 6 and a minimum number of training rows at each node of 2. Samples of the training dataset were created with 50% the size of the original dataset. This was to force some variety in the dataset subsamples used to train each tree. The default for bagging is to have the size of sample datasets match the size of the original training dataset.

A series of 4 different numbers of trees were evaluated to show the behavior of the algorithm.

The accuracy from each fold and the mean accuracy for each configuration are printed. We can see a trend of some minor lift in performance as the number of trees is increased.

## Get Your Start in Machine Learning ×

You can master applied Machine Learning **without the math or fancy degree.**

Find out how in this *free* and *practical* email course.

START MY EMAIL COURSE

Get Your Start in Machine Learning

```

1 Trees: 1
2 Scores: [87.8048780487805, 65.85365853658537, 65.85365853658537, 65.85365853658537, 73.17073170731707]
3 Mean Accuracy: 71.707%
4
5 Trees: 5
6 Scores: [60.97560975609756, 80.48780487804879, 78.04878048780488, 82.92682926829268, 63.41463414634146]
7 Mean Accuracy: 73.171%
8
9 Trees: 10
10 Scores: [60.97560975609756, 73.17073170731707, 82.92682926829268, 80.48780487804879, 68.29268292682927]
11 Mean Accuracy: 73.171%
12
13 Trees: 50
14 Scores: [63.41463414634146, 75.60975609756098, 80.48780487804879, 75.60975609756098, 85.36585365853658]
15 Mean Accuracy: 76.098%

```

A difficulty of this method is that even though deep trees are constructed, the bagged trees that are created are very similar. In turn, the predictions made by these trees are also similar, and the high variance we desire among the trees trained on different

This is because of the greedy algorithm used in the construction of the trees selecting the same or s

The tutorial tried to re-inject this variance by constraining the sample size used to train each tree. A that may be evaluated when creating each split point. This is the method used in the Random Forest

## Extensions

- **Tune the Example.** Explore different configurations for the number of trees and even individual results.
- **Bag Another Algorithm.** Other algorithms can be used with bagging. For example, a k-nearest high variance and is a good candidate for bagging.
- **Regression Problems.** Bagging can be used with regression trees. Instead of predicting the most common class value from the set of predictions, you can return the average of the predictions from the bagged trees. Experiment on regression problems.

**Did you try any of these extensions?**  
Share your experiences in the comments below.

## Review

### Get Your Start in Machine Learning

You can master applied Machine Learning **without the math or fancy degree.**

Find out how in this *free* and *practical* email course.

START MY EMAIL COURSE

In this tutorial, you discovered how to implement bootstrap aggregation from scratch with Python.

Specifically, you learned:

- How to create a subsample and estimate bootstrap quantities.
- How to create an ensemble of decision trees and use them to make predictions.
- How to apply bagging to a real world predictive modeling problem.

### Do you have any questions?

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

## Want to Code Algorithms in Python With

### Code Your First Algorithm

...with step-by-step tutorials on r

Discover how in my new  
[Machine Learning Algorithms](#)

It covers **18 tutorials** with all the code for  
Linear Regression, k-Nearest Neighbors, Stochastic

### Finally, Pull Back the

### Machine Learning Algorithms

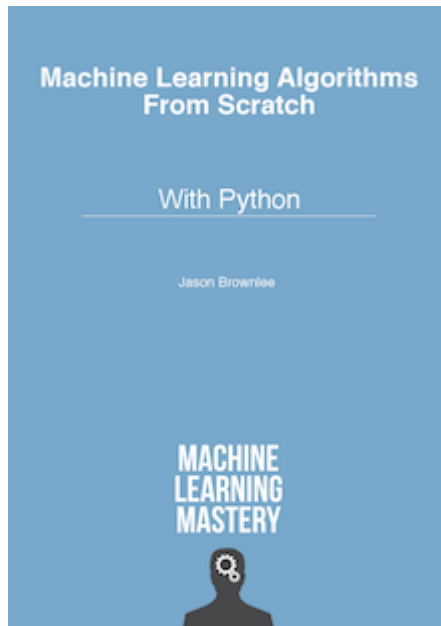
Skip the Academics. Just Results.

[Click to learn more.](#)

## Get Your Start in Machine Learning

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

START MY EMAIL COURSE





### About Jason Brownlee

Dr. Jason Brownlee is a husband, proud father, academic researcher, author, professional developer and a machine learning practitioner. He is dedicated to helping developers get started and get good at applied machine learning. [Learn more.](#)

[View all posts by Jason Brownlee](#) →

< [How To Implement The Decision Tree Algorithm From Scratch In Python](#)

[How to Implement Random Forest From Scratch in Python](#) >

## 6 Responses to *How to Implement Bagging From Scratch With Python*



**skorzec** January 18, 2017 at 4:30 am #

Thanks for this example in python from scratch. I think the reason your score stays the same is your split attributes. This leads to similar trees and thus a small variance of the ensemble set. If you change the split attributes, you will see an increase in performance of the ensemble. One solution is to alter this code s

# Build a decision tree

```
def build_tree(train, max_depth, min_size):
```

```
    root = get_split(train)
```

```
    split(root, max_depth, min_size, 1)
```

```
    return root
```

**Jason Brownlee** January 18, 2017 at 10:17 am #

### Get Your Start in Machine Learning ×

You can master applied Machine Learning **without the math or fancy degree.**

Find out how in this *free* and *practical* email course.

[START MY EMAIL COURSE](#)

[Get Your Start in Machine Learning](#)



Thanks for the tip!



**Jovan Sardinha** June 25, 2017 at 9:41 am #

REPLY ↩

thanks for the great tutorial!

I have re-written this script using sklearn for easy implementation:

<https://gist.github.com/JovanSardinha/2c58bd1e7e3aa4c02affedfe7abe8a29>



**Jason Brownlee** June 26, 2017 at 6:04 am #

Nice work!



**Alex Godfrey** August 31, 2017 at 6:43 am #

Hi Jason,

Minor tip – In the string to integer conversion – I found that the unique set gets created in some scripts that use this function. To avoid this I have changed the line to read-

```
unique = sorted(set(class_values))
```

This results in creating the same lookup dictionary every time. I ran across this when I was using the lookup dictionary to create nicely mapped printouts of the intermediate data.

## Get Your Start in Machine Learning



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

START MY EMAIL COURSE



**Jason Brownlee** September 1, 2017 at 6:38 am #

REPLY ↩

Great tip, thanks Alex!

Get Your Start in Machine Learning

## Leave a Reply

Name (required)

Email (will not be published) (required)

Website

### Get Your Start in Machine Learning ×

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

## Welcome to Machine Learning Mastery



Hi, I'm Dr. Jason Brownlee.  
My goal is to make practitioners like YOU awesome at applied machine learning.

[Read More](#)[Get Your Start in Machine Learning](#)



## Code Algorithms From Scratch in Python

Discover how to code top machine learning algorithms from first principles with Python.

Code Machine Learning Algorithms From Scratch Today!



## Get Your Start in Machine Learning

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

START MY EMAIL COURSE

### POPULAR



**Time Series Prediction with LSTM Recurrent Neural Networks in Python with Keras**

JULY 21, 2016



**Your First Machine Learning Project in Python Step-By-Step**

JUNE 10, 2016



**Develop Your First Neural Network in Python With Keras Step-By-Step**

MAY 24, 2016

Get Your Start in Machine Learning



### Sequence Classification with LSTM Recurrent Neural Networks in Python with Keras

JULY 26, 2016



### How to Setup a Python Environment for Machine Learning and Deep Learning with Anaconda

MARCH 13, 2017



### Time Series Forecasting with the Long Short-Term Memory Network in Python

APRIL 7, 2017



### Multi-Class Classification Tutorial with the Keras Deep Learning Library

JUNE 2, 2016



### Regression Tutorial with the Keras Deep Learning Library in Python

JUNE 9, 2016



### Multivariate Time Series Forecasting with LSTMs in Keras

AUGUST 14, 2017



### How to Implement the Backpropagation Algorithm From Scratch In Python

NOVEMBER 7, 2016

## Get Your Start in Machine Learning ×

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

START MY EMAIL COURSE

© 2017 Machine Learning Mastery. All Rights Reserved.

[Privacy](#) | [Contact](#) | [About](#)

Get Your Start in Machine Learning