



## BUILDING RANDOM FOREST CLASSIFIER WITH PYTHON SCIKIT LEARN

 June 26, 2017  Saimadhu Polamuri  17 Comments  Machine Learning

# Classify Breast Cancer Benign Or Malignant With Random Forest Algorithm

dataaspirant.com

Random Forest Algorithm in Python

## Building Random Forest Algorithm in Python

In the Introductory article about [random forest algorithm](#), we addressed how the random forest algorithm works with real life examples. As continues to that, In this article we are going to build the **random forest algorithm** in python with the help of one of the best [Python machine learning library Scikit-Learn](#).

To build the random forest algorithm we are going to use the Breast Cancer dataset. To summarize in this article we are going to build a random forest classifier to predict the Breast cancer type (**Benign or Malignant**).

“ Before we begin. Let’s quickly look at the **table of contents**.

### Table of contents:

- Overview of Random forest algorithm
- About Breast Cancer
  - Benign
  - Malignant
- UCI breast cancer dataset description
- Machine learning workflow
- Implementing random forest algorithm in Python
  - Creating dataset
  - Handling missing values
  - Split data into train and test dataset
  - Training random forest classifier with scikit learn
- Perform predictions
- Accuracy calculations
  - Train Accuracy
  - Test Accuracy
- Confusion matrix
- Summary
- Recommended Data Science Courses

- [Decision trees and random forest algorithm guide.](#)
- [Ensemble machine learning: Random forest and Adaboost.](#)
- [Machine learning classification concepts for beginners.](#)
- [Applying machine learning classification techniques case studies.](#)

## Building Random Forest Algorithm in Python

CLICK TO TWEET 

### Overview of Random forest algorithm

Random forest algorithm is an **ensemble classification** algorithm. Ensemble classifier means a group of classifiers. Instead of using only one classifier to predict the target, In ensemble, we use multiple classifiers to predict the target.

In case, of random forest, these ensemble classifiers are the randomly created decision trees. Each [decision tree](#) is a single classifier and the target prediction is based on the **majority voting** method.

The majority voting concept is same as the political votings. Each person **votes** per one political party out all the political parties participating in elections. In the same way, **every classifier will votes** to one target class out of all the target classes.

To declare the election results. The votes will calculate and the party which got the **most number of votes** treated as the election winner. In the same way, the target class which got the **most number of votes** considered as the final predicted target class.

Before we go further it's better to spend some time on the below articles to understand how the random forest algorithm works.

- “ How the random forest algorithm works in machine learning
- “ Decision tree algorithm concepts for beginners
- “ Building decision tree classifier in Python

I hope you have a clear understanding of how the random forest algorithm works. Now let's implement the same. As I said earlier, we are going to use the breast cancer dataset to implement the random forest.

Before we begin let's look at some stats and the **impact of breast cancer** in present generation.

## About Breast Cancer

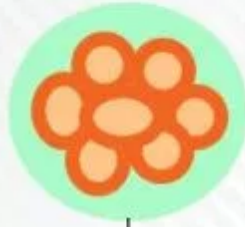
Sadly breast cancer is to **second most** death reason for women's. In the US during the year 2016, almost **246,660** women's breast cancer cases are diagnosed. The myth people believe tumor as cancer but which is not true.

Only the **continuously growing tumor** causes death. Based on this properties the tumors are mainly of 2 kinds.

- Benign Tumor
- Malignant Tumor



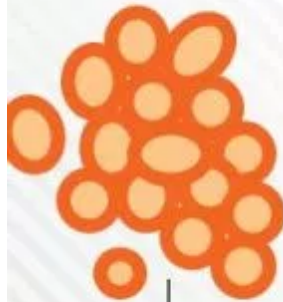
## *Learn the difference and what you can do to prevent them!*



**BENIGN** (not cancer) tumor cells grow only locally and cannot spread by invasion or metastasis.

They are often surrounded by a protective “sac”  
– a mechanism performed by your immune system  
– that segregates it from the rest of your body and enables it to be easily removed.

- CAPSULATED
- SLOW GROWING
- NON-INVASIVE
- NON-CANCEROUS



**MALIGNANT** (cancer) cells invade neighboring tissues, enter blood vessels, and metastasize to different sites.

They may not have symptoms initially and the first indication that something isn't right may be the detection of a painless lump.

- CANCEROUS
- FAST GROWING
- INVASIVE & INFILTRATE
- NON-CAPSULATED



malignant benign tumor difference  
Image Credit:: thetruthaboutcancer.com

### **Benign:**

A benign tumor is not a cancerous tumor. Which means it's not able to spread through the body like the cancerous tumors. The benign is serious when it's growing in sensitive places. This kind of tumors are will well terminated with proper treatment and with the change in diet habits.

### **Malignant**

The malignant tumor is the cancerous tumor which causes death. These tumors can grow so fast and spread over various parts of the body.

A good read about these tumor and health prevention can be found in the [thetruthaboutcancer](https://www.thetruthaboutcancer.com/) article.

### **UCI breast cancer dataset description**

We are using the UCI breast cancer dataset to build the random forest classifier in Python. You can download the data from [UCI](https://archive.ics.uci.edu/dataset/16/breast+cancer) or You can download the code from [Dataaspirant Github](https://github.com/Dataaspirant).

This breast cancer dataset is the most popular classification dataset. Which is having **10 features** and **1 target class**.

### **Breast Cancer dataset features:**

- Sample code number:
  - id number
- Clump Thickness:
  - The values are in the range of 1 – 10
- Uniformity of Cell Size:
  - The values are in the range of 1 – 10
- Uniformity of Cell Shape:
  - The values are in the range of 1 – 10
- Marginal Adhesion:
  - The values are in the range of 1 – 10
- Single Epithelial Cell Size:
  - The values are in the range of 1 – 10
- Bare Nuclei:
  - The values are in the range of 1 – 10
- Bland Chromatin:
  - The values are in the range of 1 – 10
- Normal Nucleoli:
  - The values are in the range of 1 – 10
- Mitoses:
  - The values are in the range of 1 – 10



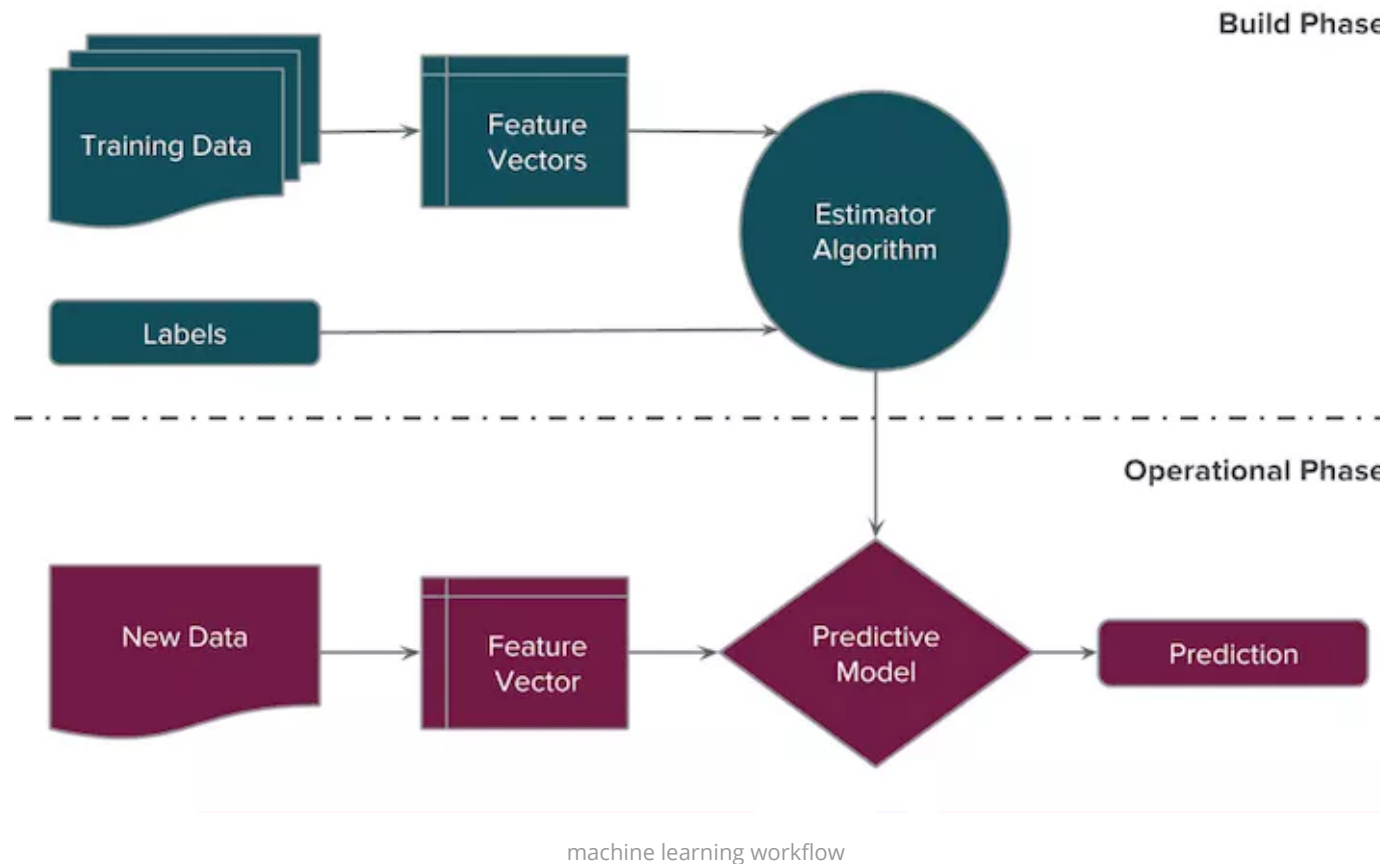
## Breast Cancer dataset Target:

The target class having two target classes

- Bening
  - The value will be 2
- Malignant
  - The value will be 4

This dataset also having **missing values**. In the coding section of this article, we are to going deal with the missing values before we model the random forest algorithm.

## Machine learning workflow



## Implementing random forest algorithm in Python

To implement the random forest algorithm we are going to follow the below **two phase** with step by step workflow.

- **Build Phase**
  - Creating dataset
  - Handling missing values

- Splitting data into train and test datasets
- Training random forest classifier with Python scikit learn
- **Operational Phase**
  - Perform predictions
  - Accuracy calculations
    - Train Accuracy
    - Test Accuracy

Let's begin the journey of building the random forest classifier with importing the required Python machine learning packages.

### Import required Python machine learning packages

```
1 # Required Python Packages
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import accuracy_score
6 from sklearn.metrics import confusion_matrix
```

The above python machine learning packages we are going to use to build the random forest classifier. Let's talk about the need for these packages in random forest classifier implementation.

- **Pandas:**
  - Pandas package is the best choice for tabular data analysis.
  - All the data manipulation tasks in this article are going to use the Pandas methods.
- **train\_test\_split:**
  - We imported scikit-learn **train\_test\_split** method to split the breast cancer dataset into test and train dataset.
  - Train dataset will be used in the training phase and the test dataset will be used in the validation phase.
- **RandomForestClassifier:**

- We imported scikit-learn **RandomForestClassifier** method to model the training dataset with random forest classifier.
- Later the modeled random forest classifier used to perform the predictions.
- **accuracy\_score:**
  - We imported scikit-learn **accuracy\_score** method to calculate the accuracy of the trained classifier.
- **confusion\_matrix:**
  - We imported scikit-learn **confusion\_matrix** to understand the trained classifier behavior over the test dataset or validate dataset.

Copy the above code in any text file (or you favorite txt editor) and save the file with the python extension (**.py**). Let say **random\_forest.py**

Then call the random\_forest.py file from the terminal using the below command.

```
1 python random_forest.py
```

If you install the python [machine learning packages](#) properly, you won't face any issues. Even though you install the packages properly and you facing the issue **ImportError: No module named model\_selection**. This means the scikit learn package you are using not updated to the new version.

I hope you are using scikit learn **0.17 or lesser version**. You can copy and paste the below code to know your scikit learn version.

```
1 import sklearn
2 print (sklearn.__version__)
```

If the version your are using is **0.17** or lesser than that, you need to update your scikit learn version to **0.18**

You can use the below commands to update your scikit learn to the new version (0.18)

## Using Pip

```
1 pip install -U scikit-learn
```

## Using Anaconda

```
1 conda install scikit-learn=0.18
```

Once you upgraded your scikit-learn package. Run the above code and you won't face any issues. If you still face any issue to run the above code do please let me know in the **comments section**.

Now let's create the dataset to model the random forest classifier.

## Creating dataset

The downloaded dataset is in the **data** format. So we are going to convert into the **CSV** format. To do that we are going to write a simple function which first loads the **data** format into the pandas dataframe and later the loaded dataframe will save into the **CSV** file format.

```
1 # File Paths
2 INPUT_PATH = "../inputs/breast-cancer-wisconsin.data"
3 OUTPUT_PATH = "../inputs/breast-cancer-wisconsin.csv"
4
5
6 def data_file_to_csv():
7     """
8
9     :return:
10    """
11
12    # Headers
13    headers = ["CodeNumber", "ClumpThickness", "UniformityCellSize", "UniformityCellShape", "MarginalAdhesion",
14               "SingleEpithelialCellSize", "BareNuclei", "BlandChromatin", "NormalNucleoli", "Mitoses",
15               "CancerType"]
16
17    # Load the dataset into Pandas data frame
18    dataset = read_data(INPUT_PATH)
```

```
18 # Add the headers to the loaded dataset
19 dataset = add_headers(dataset, headers)
20 # Save the loaded dataset into csv format
21 dataset.to_csv(OUTPUT_PATH, index=False)
22 print "File saved ...!"
```

The **INPUT\_PATH** is having the path for the downloaded data format file and the **OUTPUT\_PATH** is having the output where the CSV format file is going to save.

Using the pandas **read\_csv** method we loaded the data format file into pandas dataframe.

The loaded dataset doesn't have the header names. So we need to add the header names to the loaded dataframe. To do the same we have written a function with takes the **dataset** and **header names** as input and add the header names to the dataset.

```
1 def add_headers(dataset, headers):
2     """
3     Add the headers to the dataset
4     :param dataset:
5     :param headers:
6     :return:
7     """
8     dataset.columns = headers
9     return dataset
```

After adding the header names to dataset we are saving the dataset into CSV format. While saving the file we parameterized the **index=False**. When we save the loaded dataframe without this the saved file will have an extra column with the indexes. So to eliminate this we are parameterized the **index=False**.

Now we are ready with the dataset. The next biggest thing is the **preprocessing** the data.

Sometimes, if it's our day, we don't need to do much work on the preprocessing stage. But, if not.

“ We need to spend a lot of time in the preprocessing stage.

Handling of missing values once such task in preprocessing the data.

## Handling missing values

The process of handling missing values will differ from dataset to dataset. For the cancer dataset, we are using simple tasks to handle the missing values in the loaded dataset.

Before reviling what those missing values, I want to show the ways to identify the missing values. So when you are working on a different dataset. You can identify the missing values by yourself.

The best idea to start with is, calculating basic statistics for each column (features and target) of the dataset. You may be wondering what the use of calculating basic statistics of the dataset and how it gonna helps to find the missing values.

**Yes**, finding the basic statistics will helps us to find the missing values in the dataset. The idea is we can use pandas describe method on the loaded dataset to calculate the basic statistics. This outputs the stats about only the columns which are not having any missing values or categorical values.

**Feeling missed somewhere**, No issue let's implement a function to calculate the basic statistics then you get the clear idea of what I talking about.

```
1 def dataset_statistics(dataset):
2     """
3     Basic statistics of the dataset
4     :param dataset: Pandas dataframe
5     :return: None, print the basic statistics of the dataset
6     """
7     print dataset.describe()
```

As I said before, We are using pandas **describe** method to get the basic statistics of the dataset.

Now let's call this function and check the what it's outputting.

```
1 def main():
2     """
3     Main function
4     :return:
5     """
6     # Load the csv file into pandas dataframe
7     dataset = pd.read_csv(OUTPUT_PATH)
8     # Get basic statistics of the loaded dataset
9     dataset_statistics(dataset)
10
11 if __name__ == "__main__":
12     main()
```

Script Output

```
1 CodeNumber ClumpThickness UniformityCellSize UniformityCellShape
2 count      6.980000e+02      698.000000      698.000000      698.000000
3 mean      1.071807e+06      4.416905      3.137536      3.210602
4 std       6.175323e+05      2.817673      3.052575      2.972867
5 min       6.163400e+04      1.000000      1.000000      1.000000
6 25%       8.702582e+05      2.000000      1.000000      1.000000
7 50%       1.171710e+06      4.000000      1.000000      1.000000
8 75%       1.238354e+06      6.000000      5.000000      5.000000
9 max       1.345435e+07     10.000000     10.000000     10.000000
10
11      MarginalAdhesion SingleEpithelialCellSize BlandChromatin
12 count      698.000000      698.000000      698.000000
13 mean       2.809456      3.217765      3.438395
14 std       2.856606      2.215408      2.440056
15 min       1.000000      1.000000      1.000000
16 25%       1.000000      2.000000      2.000000
17 50%       1.000000      2.000000      3.000000
18 75%       4.000000      4.000000      5.000000
19 max      10.000000     10.000000     10.000000
20
21      NormalNucleoli Mitoses CancerType
22 count      698.000000     698.000000     698.000000
23 mean       2.869628      1.590258      2.690544
```



24	std	3.055004	1.716162	0.951596
25	min	1.000000	1.000000	2.000000
26	25%	1.000000	1.000000	2.000000
27	50%	1.000000	1.000000	2.000000
28	75%	4.000000	1.000000	4.000000
29	max	10.000000	10.000000	4.000000

If you observe the above statistics clearly you can identify that we are missing one column details. The column header we are missing is **BareNuclei**.

Now open the CSV file and check the column details of Bare Nuclei. You will find the missing values replaced with ? Which means we need to treat those values as missing values.

As we know the missing values column and the character used to represent the missing values. Now let's write a simple function will take the dataset, header\_name and missing value representing the character as input handles the missing values.

```

1 def handel_missing_values(dataset, missing_values_header, missing_label):
2     """
3     Filter missing values from the dataset
4     :param dataset:
5     :param missing_values_header:
6     :param missing_label:
7     :return:
8     """
9
10    return dataset[dataset[missing_values_header] != missing_label]
```

Now let's call this function inside the main function

```

1 def main():
2     """
3     Main function
4     :return:
5     """
6     # Headers
7     headers = ["CodeNumber", "ClumpThickness", "UniformityCellSize", "UniformityCellShape", "MarginalAdhesion",
```

```

8         "SingleEpithelialCellSize", "BareNuclei", "BlandChromatin", "NormalNucleoli", "Mitoses",
9         "CancerType"]
10     # Load the csv file into pandas dataframe
11     dataset = pd.read_csv(OUTPUT_PATH)
12     # Get basic statistics of the loaded dataset
13     dataset_statistics(dataset)
14
15     # Filter missing values
16     dataset = handel_missing_values(dataset, HEADERS[6], '?')
17 if __name__ == "__main__":
18     main()

```

Now our dataset is **missing values free**. Now let's split the data into train and test dataset. The training dataset will use to train the random forest classifier and the test dataset used the validate the model random forest classifier.

### Split data into train and test datasets

To split the data into train and test dataset, Let's write a function which takes the dataset, train percentage, feature header names and target header name as inputs and returns the **train\_x, test\_x, train\_y and test\_y** as outputs.

```

1 def split_dataset(dataset, train_percentage, feature_headers, target_header):
2     """
3     Split the dataset with train_percentage
4     :param dataset:
5     :param train_percentage:
6     :param feature_headers:
7     :param target_header:
8     :return: train_x, test_x, train_y, test_y
9     """
10
11     # Split dataset into train and test dataset
12     train_x, test_x, train_y, test_y = train_test_split(dataset[feature_headers], dataset[target_header],
13                                                         train_size=train_percentage)
14     return train_x, test_x, train_y, test_y

```

Now let's call the above function inside the **main function**.

```

1 def main():
2     """

```

```

3 Main function
4 :return:
5 """
6 # Load the csv file into pandas dataframe
7 dataset = pd.read_csv(OUTPUT_PATH)
8 # Get basic statistics of the loaded dataset
9 dataset_statistics(dataset)
10
11 # Filter missing values
12 dataset = handel_missing_values(dataset, HEADERS[6], '?')
13 train_x, test_x, train_y, test_y = split_dataset(dataset, 0.7, HEADERS[1:-1], HEADERS[-1])
14
15
16 if __name__ == "__main__":
17     main()

```

Where the **HEADERS[1:-1]** contains all the features header names. We eliminated the first and last header names. The first header name is **id** and the last header the **target** header. The **HEADERS[-1]** contains the target header name.

We can print the shape of the **train\_x**, **test\_x**, **train\_y** and **test\_y** to check whether the split proper or not.

```

1 def main():
2     """
3     Main function
4     :return:
5     """
6     # Load the csv file into pandas dataframe
7     dataset = pd.read_csv(OUTPUT_PATH)
8     # Get basic statistics of the loaded dataset
9     dataset_statistics(dataset)
10
11     # Filter missing values
12     dataset = handel_missing_values(dataset, HEADERS[6], '?')
13     train_x, test_x, train_y, test_y = split_dataset(dataset, 0.7, HEADERS[1:-1], HEADERS[-1])
14
15     # Train and Test dataset size details
16     print "Train_x Shape :: ", train_x.shape
17     print "Train_y Shape :: ", train_y.shape
18     print "Test_x Shape :: ", test_x.shape
19     print "Test_y Shape :: ", test_y.shape
20
21 if __name__ == "__main__":

```

### Script output:

```
1 Train_x Shape :: (477, 9)
2 Train_y Shape :: (477,)
3 Test_x Shape :: (205, 9)
4 Test_y Shape :: (205,)
```

From the above result, it's clear that the train and test split was proper. Now let's build the random forest classifier using the **train\_x** and **train\_y** datasets.

### Training random forest classifier with scikit learn

To train the random forest classifier we are going to use the below **random\_forest\_classifier** function. Which requires the features (**train\_x**) and target (**train\_y**) data as inputs and returns the train random forest classifier as output.

```
1 def random_forest_classifier(features, target):
2     """
3     To train the random forest classifier with features and target data
4     :param features:
5     :param target:
6     :return: trained random forest classifier
7     """
8     clf = RandomForestClassifier()
9     clf.fit(features, target)
10    return clf
```

Now let's call the above function inside the **main function** and print the trained classifier.

```
1 def main():
2     """
3     Main function
4     :return:
5     """
```

```

6 # Load the csv file into pandas dataframe
7 dataset = pd.read_csv(OUTPUT_PATH)
8
9 # Filter missing values
10 dataset = handel_missing_values(dataset, HEADERS[6], '?')
11 train_x, test_x, train_y, test_y = split_dataset(dataset, 0.7, HEADERS[1:-1], HEADERS[-1])
12
13 # Create random forest classifier instance
14 trained_model = random_forest_classifier(train_x, train_y)
15 print "Trained model :: ", trained_model
16 if __name__ == "__main__":
17     main()

```

### Script Output:

```

1 Trained model :: RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
2     max_depth=None, max_features='auto', max_leaf_nodes=None,
3     min_impurity_split=1e-07, min_samples_leaf=1,
4     min_samples_split=2, min_weight_fraction_leaf=0.0,
5     n_estimators=10, n_jobs=1, oob_score=False, random_state=None,
6     verbose=0, warm_start=False)

```

## Perform predictions

As we model the classifier. Now let's take few observations from test dataset and print what our model is predicting and what the actual target.

To do that first let's predict target for all the test features (**test\_x**) using the trained classifier. Later we will see what our trained model is predicting and what the actual output could be.

```

1 def main():
2     """
3     Main function

```

```

4     :return:
5     """
6     # Load the csv file into pandas dataframe
7     dataset = pd.read_csv(OUTPUT_PATH)
8
9     # Filter missing values
10    dataset = handel_missing_values(dataset, HEADERS[6], '?')
11    train_x, test_x, train_y, test_y = split_dataset(dataset, 0.7, HEADERS[1:-1], HEADERS[-1])
12
13    # Create random forest classifier instance
14    trained_model = random_forest_classifier(train_x, train_y)
15    print "Trained model :: ", trained_model
16    predictions = trained_model.predict(test_x)
17
18    for i in xrange(0, 5):
19        print "Actual outcome :: {} and Predicted outcome :: {}".format(list(test_y)[i], predictions[i])
20 if __name__ == "__main__":
21     main()

```

First I converted the test\_y into list object from pandas dataframe. The reason is as we randomly split the train and test dataset the **indexes** of the test\_y won't be **in order**. If we convert the dataframe in to list object the indexes will be in order.

From the above code, we are printing the first 5 values of test\_y and the predict results target.

```

1 Actual outcome :: 2 and Predicted outcome :: 2
2 Actual outcome :: 2 and Predicted outcome :: 2
3 Actual outcome :: 2 and Predicted outcome :: 2
4 Actual outcome :: 2 and Predicted outcome :: 2
5 Actual outcome :: 4 and Predicted outcome :: 4

```

Seems like the trained classifier predicted the first 5 target classes correctly. To know more about the model let's check the train and test accuracy information.

## Accuracy calculations

```

1 def main():

```

```

2  """
3  Main function
4  :return:
5  """
6  # Load the csv file into pandas dataframe
7  dataset = pd.read_csv(OUTPUT_PATH)
8
9  # Filter missing values
10 dataset = handel_missing_values(dataset, HEADERS[6], '?')
11 train_x, test_x, train_y, test_y = split_dataset(dataset, 0.7, HEADERS[1:-1], HEADERS[-1])
12
13 # Create random forest classifier instance
14 trained_model = random_forest_classifier(train_x, train_y)
15 print "Trained model :: ", trained_model
16 predictions = trained_model.predict(test_x)
17
18 # Train and Test Accuracy
19 print "Train Accuracy :: ", accuracy_score(train_y, trained_model.predict(train_x))
20 print "Test Accuracy :: ", accuracy_score(test_y, predictions)
21
22 if __name__ == "__main__":
23     main()

```

To calculate the accuracy we are using scikit learn the **accuracy\_score** method.

### Script output:

```

1 Train Accuracy :: 0.991614255765
2 Test Accuracy  :: 0.970731707317

```

Our trained classifier model giving **99%** accuracy for train dataset and **97%** accuracy for test dataset.

### Confusion matrix

To know the ture\_postive and true\_negative details. Let's print the confusion matrix of our trained classifier.

```

1 def main():
2     """
3     Main function

```

```

4     :return:
5     """
6     # Load the csv file into pandas dataframe
7     dataset = pd.read_csv(OUTPUT_PATH)
8
9     # Filter missing values
10    dataset = handel_missing_values(dataset, HEADERS[6], '?')
11    train_x, test_x, train_y, test_y = split_dataset(dataset, 0.7, HEADERS[1:-1], HEADERS[-1])
12
13    # Create random forest classifier instance
14    trained_model = random_forest_classifier(train_x, train_y)
15    predictions = trained_model.predict(test_x)
16
17    print "Train Accuracy :: ", accuracy_score(train_y, trained_model.predict(train_x))
18    print "Test Accuracy :: ", accuracy_score(test_y, predictions)
19    print " Confusion matrix ", confusion_matrix(test_y, predictions)
20
21
22 if __name__ == "__main__":
23     main()

```

### Script output:

```

1 Train Accuracy :: 0.991614255765
2 Test Accuracy  :: 0.970731707317
3 Confusion matrix [[123  5]
4 [ 1 76]]

```

You can get the complete code below. If you know GitHub, you can clone the complete code from our [Github](#) account.

### Complete Code:

```

1 # Required Python Packages
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import accuracy_score
6 from sklearn.metrics import confusion_matrix

```



```

7
8 import pdb
9
10 # File Paths
11 INPUT_PATH = "../inputs/breast-cancer-wisconsin.data"
12 OUTPUT_PATH = "../inputs/breast-cancer-wisconsin.csv"
13
14 # Headers
15 HEADERS = ["CodeNumber", "ClumpThickness", "UniformityCellSize", "UniformityCellShape", "MarginalAdhesion",
16            "SingleEpithelialCellSize", "BareNuclei", "BlandChromatin", "NormalNucleoli", "Mitoses", "CancerType"]
17
18
19 def read_data(path):
20     """
21     Read the data into pandas dataframe
22     :param path:
23     :return:
24     """
25     data = pd.read_csv(path)
26     return data
27
28
29 def get_headers(dataset):
30     """
31     dataset headers
32     :param dataset:
33     :return:
34     """
35     return dataset.columns.values
36
37
38 def add_headers(dataset, headers):
39     """
40     Add the headers to the dataset
41     :param dataset:
42     :param headers:
43     :return:
44     """
45     dataset.columns = headers
46     return dataset
47
48
49 def data_file_to_csv():
50     """
51

```

```

52     :return:
53     """
54
55     # Headers
56     headers = ["CodeNumber", "ClumpThickness", "UniformityCellSize", "UniformityCellShape", "MarginalAdhesion",
57               "SingleEpithelialCellSize", "BareNuclei", "BlandChromatin", "NormalNucleoli", "Mitoses",
58               "CancerType"]
59     # Load the dataset into Pandas data frame
60     dataset = read_data(INPUT_PATH)
61     # Add the headers to the loaded dataset
62     dataset = add_headers(dataset, headers)
63     # Save the loaded dataset into csv format
64     dataset.to_csv(OUTPUT_PATH, index=False)
65     print "File saved ...!"
66
67
68 def split_dataset(dataset, train_percentage, feature_headers, target_header):
69     """
70     Split the dataset with train_percentage
71     :param dataset:
72     :param train_percentage:
73     :param feature_headers:
74     :param target_header:
75     :return: train_x, test_x, train_y, test_y
76     """
77
78     # Split dataset into train and test dataset
79     train_x, test_x, train_y, test_y = train_test_split(dataset[feature_headers], dataset[target_header],
80                                                         train_size=train_percentage)
81     return train_x, test_x, train_y, test_y
82
83
84 def handel_missing_values(dataset, missing_values_header, missing_label):
85     """
86     Filter missing values from the dataset
87     :param dataset:
88     :param missing_values_header:
89     :param missing_label:
90     :return:
91     """
92
93     return dataset[dataset[missing_values_header] != missing_label]
94
95
96 def random_forest_classifier(features, target):

```

```

97     """
98     To train the random forest classifier with features and target data
99     :param features:
100     :param target:
101     :return: trained random forest classifier
102     """
103     clf = RandomForestClassifier()
104     clf.fit(features, target)
105     return clf
106
107
108 def dataset_statistics(dataset):
109     """
110     Basic statistics of the dataset
111     :param dataset: Pandas dataframe
112     :return: None, print the basic statistics of the dataset
113     """
114     print dataset.describe()
115
116
117 def main():
118     """
119     Main function
120     :return:
121     """
122     # Load the csv file into pandas dataframe
123     dataset = pd.read_csv(OUTPUT_PATH)
124     # Get basic statistics of the loaded dataset
125     dataset_statistics(dataset)
126
127     # Filter missing values
128     dataset = handel_missing_values(dataset, HEADERS[6], '?')
129     train_x, test_x, train_y, test_y = split_dataset(dataset, 0.7, HEADERS[1:-1], HEADERS[-1])
130
131     # Train and Test dataset size details
132     print "Train_x Shape :: ", train_x.shape
133     print "Train_y Shape :: ", train_y.shape
134     print "Test_x Shape :: ", test_x.shape
135     print "Test_y Shape :: ", test_y.shape
136
137     # Create random forest classifier instance
138     trained_model = random_forest_classifier(train_x, train_y)
139     print "Trained model :: ", trained_model
140     predictions = trained_model.predict(test_x)
141

```

```
142     for i in xrange(0, 5):
143         print "Actual outcome :: {} and Predicted outcome :: {}".format(list(test_y)[i], predictions[i])
144
145     print "Train Accuracy :: ", accuracy_score(train_y, trained_model.predict(train_x))
146     print "Test Accuracy :: ", accuracy_score(test_y, predictions)
147     print " Confusion matrix ", confusion_matrix(test_y, predictions)
148
149
150 if __name__ == "__main__":
151     main()
```

## Summary

In this article, you learned how to implement the most popular classification algorithm random forest in python using python scikit learn package.

On process, you learned how to handle the missing values. Finally, you learned how to calculate the accuracy of any trained classifier using the scikit learn **accuracy\_score** method.

### Follow us:

[FACEBOOK](#) | [QUORA](#) | [TWITTER](#) | [GOOGLE+](#) | [LINKEDIN](#) | [REDDIT](#) | [FLIPBOARD](#) | [MEDIUM](#) | [GITHUB](#)

I hope you like this post. If you have any questions, then feel free to comment below. If you want me to write on one particular topic, then do tell it to me in the comments below.

## Recommended Data Science Courses

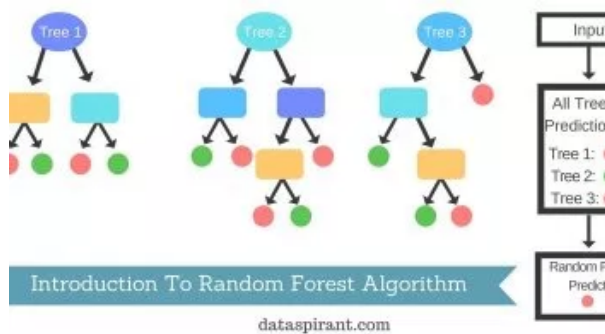
- [Decision trees and random forest algorithm guide.](#)
- [Ensemble machine learning: Random forest and Adaboost.](#)

- Machine learning classification concepts for beginners.
- Applying machine learning classification techniques case studies.

Share this:



## Related



How the random forest algorithm works in machine learning

May 22, 2017

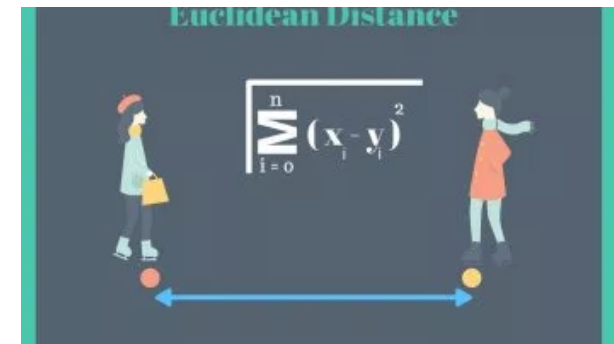
In "Machine Learning"



visualize decision tree in python with graphviz

April 21, 2017

In "Machine Learning"



Knn sklearn, K-Nearest Neighbor implementation with scikit learn

December 30, 2016

In "Data Science"

classification algorithms

python

random forest

## Previous Post

< [How the random forest algorithm works in machine learning](#)

## Next Post

[How to perform the principal component analysis in R](#) >

Machine Learning



Machine Learning

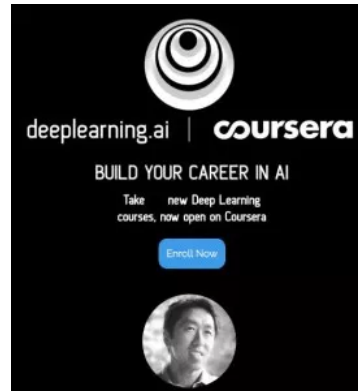
**AWARDED TOP 75  
DATA SCIENCE BLOG**



**FOLLOW AUTHOR**



## BUILD YOUR CAREER IN AI WITH ANDREW NG DEEP LEARNING COURSES



Andrew ng Deep learning courses

## MOST POPULAR POSTS



How Decision Tree Algorithm works



Building Decision Tree Algorithm in Python with scikit learn



Difference Between Softmax Function and Sigmoid Function



How the random forest algorithm works in machine learning



Five most popular similarity measures implementation in python



Support Vector  
Machine Classifier  
Implementation in  
R with caret  
package



Building Random  
Forest Classifier  
with Python Scikit  
learn



visualize decision  
tree in python  
with graphviz

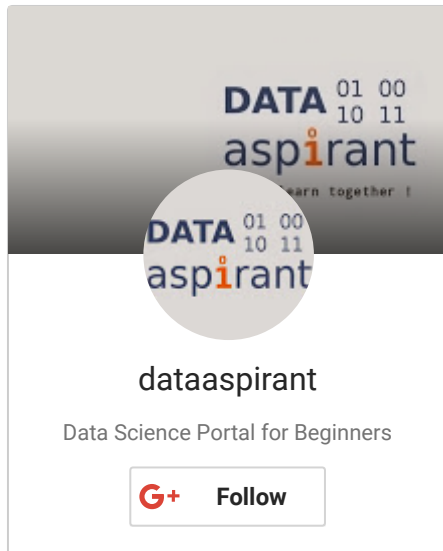


How the Naive  
Bayes Classifier  
works in Machine  
Learning



Support vector  
machine (Svm  
classifier)  
implemenation in  
python with Scikit-  
learn





## AEGIS GRAHAM BELL AWARDS 2018



## RECENT POSTS

Twitter Sentiment analysis using R

How Q learning can be used in reinforcement learning

How to perform Reinforcement learning with R

Feature selection techniques with R

How to perform hierarchical clustering in R

10 Smart R programming tips to become better R programmer

How to create histograms in R

How to perform the principal component analysis in R

Building Random Forest Classifier with Python Scikit learn

How the random forest algorithm works in machine learning



## FOLLOW ON TWITTER

---

## CATEGORIES

---

Amazon

Big Data

Courses

Data Science

Data Science Events

DATAMINING

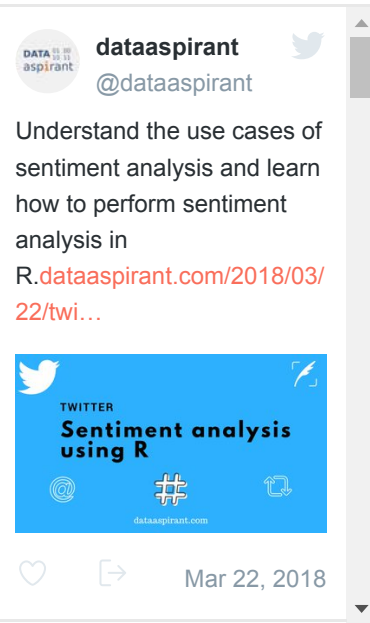
Deep Learning

## NEVER MISS A BIT

---

Hey Dude Subscribe to Dataaspirant. To get post updates in your inbox.

## Tweets by @dataaspirant



[Embed](#)

[View on Twitter](#)

Interviews

Machine Learning

Newsletter

python

R

Recommendation  
Engine

### QUICK LINKS

[Home](#)  
[About](#)  
[Join us](#)  
[Beginners Guide](#)  
[Data scientists Interviews](#)

### HOURS & INFO

WhiteField, Bangalore,  
India  
Opens: 9:00 AM  
Closes: 11:00 PM

### CATEGORIES

[Amazon](#)  
[Big Data](#)  
[Courses](#)  
[Data Science](#)  
[Data Science Events](#)

### TAGS

[Amazon go](#) [Big data](#) [Bigdata](#)  
[Classification](#)  
[classification](#)  
[algorithms](#) [clustering](#)

partners

Data science courses

DATAMINING

Deep Learning

Interviews

Machine Learning

Newsletter

python

R

Recommendation Engine

algorithms datamining Data  
mining DataScience data  
science DataScienceCongress2017  
Data science Courses Data Science  
Events data scientist Decision  
tree deep learning hierarchical  
clustering k-nearest neighbor  
kaggle Linear Regression  
logistic regression  
Machine learning  
monthly newsletter  
multinomial logistic regression  
Naive Bayes pca python  
python programming  
language random forest  
Recommendation\_engine  
Recommendation\_systems  
Regression regression coefficient  
Reinforcement Learning R  
Programming Scala  
scikit-learn SciPy Script output  
similarity\_distance Supervised  
Learning Svm tensorflow Unsupervised  
Learning virtual\_env

