

[Log in](#)[Create Account](#)

Avinash Navlani
December 28th, 2018

PYTHON

Decision Tree Classification in Python

In this tutorial, learn Decision Tree Classification, attribute selection measures, and how to build and optimize Decision Tree Classifier using Python Scikit-learn package.

As a marketing manager, you want a set of customers who are most likely to purchase your product. This is how you can save your marketing budget by finding your audience. As a loan manager, you need to identify risky loan applications to achieve a lower loan default rate. This process of classifying customers into a group of potential and non-potential customers or safe or risky loan applications is known as a classification problem. Classification is a two-step process, learning step and prediction step. In the learning step, the model is developed based on given training data. In the prediction step, the model is used to predict the response for given data. Decision Tree is one of the easiest and popular classification algorithms to understand and interpret. It can be utilized for both classification and regression kind of problem.

In this tutorial, you are going to cover the following topics:

- Decision Tree Algorithm
[Want to leave a comment?](#)

• How does the Decision Tree algorithm work?

- Attribute Selection Measures

 - Information Gain

 - Gain Ratio

 - Gini index

- Optimizing Decision Tree Performance

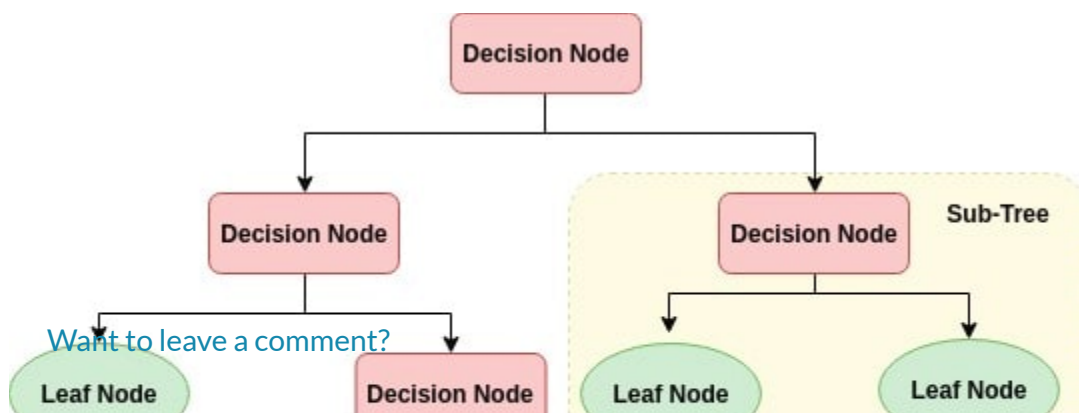
- Classifier Building in Scikit-learn

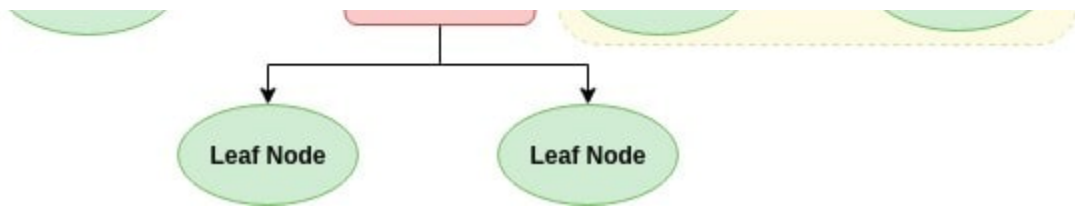
- Pros and Cons

- Conclusion

Decision Tree Algorithm

A decision tree is a flowchart-like tree structure where an internal node represents feature(or attribute), the branch represents a decision rule, and each leaf node represents the outcome. The topmost node in a decision tree is known as the root node. It learns to partition on the basis of the attribute value. It partitions the tree in recursively manner call recursive partitioning. This flowchart-like structure helps you in decision making. It's visualization like a flowchart diagram which easily mimics the human level thinking. That is why decision trees are easy to understand and interpret.



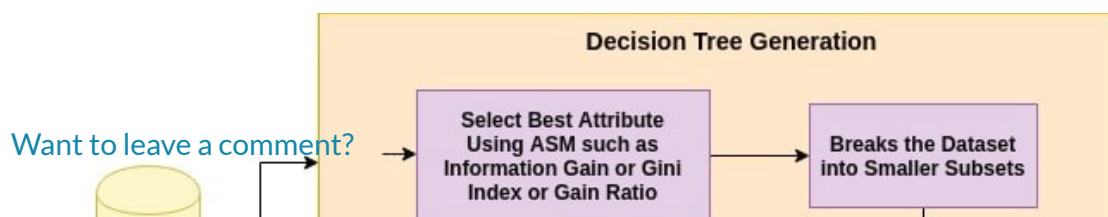


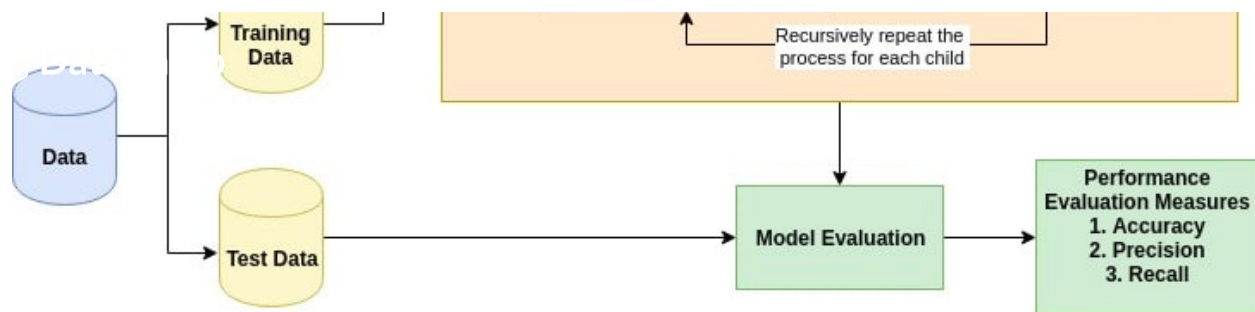
Decision Tree is a white box type of ML algorithm. It shares internal decision-making logic, which is not available in the black box type of algorithms such as Neural Network. Its training time is faster compared to the neural network algorithm. The time complexity of decision trees is a function of the number of records and number of attributes in the given data. The decision tree is a distribution-free or non-parametric method, which does not depend upon probability distribution assumptions. Decision trees can handle high dimensional data with good accuracy.

How does the Decision Tree algorithm work?

The basic idea behind any decision tree algorithm is as follows:

1. Select the best attribute using Attribute Selection Measures(ASM) to split the records.
2. Make that attribute a decision node and breaks the dataset into smaller subsets.
3. Starts tree building by repeating this process recursively for each child until one of the condition will match:
 - All the tuples belong to the same attribute value.
 - There are no more remaining attributes.
 - There are no more instances.





Attribute Selection Measures

Attribute selection measure is a heuristic for selecting the splitting criterion that partition data into the best possible manner. It is also known as splitting rules because it helps us to determine breakpoints for tuples on a given node. ASM provides a rank to each feature(or attribute) by explaining the given dataset. Best score attribute will be selected as a splitting attribute ([Source](#)). In the case of a continuous-valued attribute, split points for branches also need to define. Most popular selection measures are Information Gain, Gain Ratio, and Gini Index.

Information Gain

Shannon invented the concept of entropy, which measures the impurity of the input set. In physics and mathematics, entropy referred as the randomness or the impurity in the system. In information theory, it refers to the impurity in a group of examples. Information gain is the decrease in entropy. Information gain computes the difference between entropy before split and average entropy after split of the dataset based on given attribute values. ID3 (Iterative Dichotomiser) decision tree algorithm uses information gain.

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2 p_i$$

Where, P_i is the probability that an arbitrary tuple in D belongs to class C_i .

[Want to leave a comment?](#)

$$\text{Info}_A(D) = \sum_{j=1}^V \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

Where,

- $\text{Info}(D)$ is the average amount of information needed to identify the class label of a tuple in D .
- $|D_j|/|D|$ acts as the weight of the j th partition.
- $\text{Info}_A(D)$ is the expected information required to classify a tuple from D based on the partitioning by A .

The attribute A with the highest information gain, $\text{Gain}(A)$, is chosen as the splitting attribute at node $N()$.

Gain Ratio

Information gain is biased for the attribute with many outcomes. It means it prefers the attribute with a large number of distinct values. For instance, consider an attribute with a unique identifier such as `customer_ID` has zero $\text{info}(D)$ because of pure partition. This maximizes the information gain and creates useless partitioning.

C4.5, an improvement of ID3, uses an extension to information gain known as the gain ratio. Gain ratio handles the issue of bias by normalizing the information gain using Split Info. Java implementation of the C4.5 algorithm is known as J48, which is available in WEKA data mining tool.

$$\text{SplitInfo}_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left(\frac{|D_j|}{|D|} \right)$$

[Want to leave a comment?](#)

Where,

- $|D_j|/|D|$ acts as the weight of the j th partition.
- v is the number of discrete values in attribute A .

The gain ratio can be defined as

$$\text{GainRatio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}_A(D)}$$

The attribute with the highest gain ratio is chosen as the splitting attribute ([Source](#)).

Gini index

Another decision tree algorithm CART (Classification and Regression Tree) uses the Gini method to create split points.

$$\text{Gini}(D) = 1 - \sum_{i=1}^m p_i^2$$

Where, p_i is the probability that a tuple in D belongs to class C_i .

The Gini Index considers a binary split for each attribute. You can compute a weighted sum of the impurity of each partition. If a binary split on attribute A partitions data D into D_1 and D_2 , the Gini index of D is:

$$\text{Gini}_A(D) = \frac{|D_1|}{|D|} \text{Gini}(D_1) + \frac{|D_2|}{|D|} \text{Gini}(D_2)$$

In case of a discrete-valued attribute, the subset that gives the minimum gini index [Want to leave a comment?](#) for that chosen is selected as a splitting attribute. In the case of continuous-valued

attributes, the strategy is to select each pair of adjacent values as a possible split-point and point with smaller gini index chosen as the splitting point.

$$\Delta Gini(A) = Gini(D) - Gini_A(D).$$

The attribute with minimum Gini index is chosen as the splitting attribute.

Decision Tree Classifier Building in Scikit-learn

Importing Required Libraries

Let's first load the required libraries.

```
# Load libraries

import pandas as pd

from sklearn.tree import DecisionTreeClassifier # Import Decision Tree Classifier
from sklearn.model_selection import train_test_split # Import train_test_split functi
from sklearn import metrics #Import scikit-learn metrics module for accuracy calculat
```

Loading Data

Let's first load the required Pima Indian Diabetes dataset using pandas' read CSV function. You can download the data [here](#).

```
col_names = ['pregnant', 'glucose', 'bp', 'skin', 'insulin', 'bmi', 'pedigree', 'age']
# load dataset
pima = pd.read_csv("pima-indians-diabetes.csv", header=None, names=col_names)
```

```
pima.head()
```

Want to leave a comment?									
pregnant	glucose	bp	skin	insulin	bmi	pedigree	age	lat	

	pregnant	glucose	bp	skin	insulin	bmi	pedigree	age	label
0	6	148	72	35	0	33.6	0.627	50	1
1	1	85	66	29	0	26.6	0.351	31	0
2	8	183	64	0	0	23.3	0.672	32	1
3	1	89	66	23	94	28.1	0.167	21	0
4	0	137	40	35	168	43.1	2.288	33	1

Feature Selection

Here, you need to divide given columns into two types of variables dependent(or target variable) and independent variable(or feature variables).

```
#split dataset in features and target variable
feature_cols = ['pregnant', 'insulin', 'bmi', 'age', 'glucose', 'bp', 'pedigree']
X = pima[feature_cols] # Features
y = pima.label # Target variable
```

Splitting Data

To understand model performance, dividing the dataset into a training set and a test set is a good strategy.

Let's split the dataset by using function `train_test_split()`. You need to pass 3 parameters features, target, and test_set size.

```
# Split dataset into training set and test set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state
```

[Want to leave a comment?](#)

Building Decision Tree Model

Let's create a Decision Tree Model using Scikit-learn.

```
# Create Decision Tree classifier object
clf = DecisionTreeClassifier()

# Train Decision Tree Classifier
clf = clf.fit(X_train,y_train)

#Predict the response for test dataset
y_pred = clf.predict(X_test)
```

Evaluating Model

Let's estimate, how accurately the classifier or model can predict the type of cultivars.

Accuracy can be computed by comparing actual test set values and predicted values.

```
# Model Accuracy, how often is the classifier correct?
print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
```

```
Accuracy: 0.6753246753246753
```

Well, you got a classification rate of 67.53%, considered as good accuracy. You can improve this accuracy by tuning the parameters in the Decision Tree Algorithm.

Visualizing Decision Trees

You can use Scikit-learn's `export_graphviz` function for display the tree within a Jupyter Notebook. [Want to leave a comment?](#) For plotting tree, you also need to install graphviz and pydotplus.

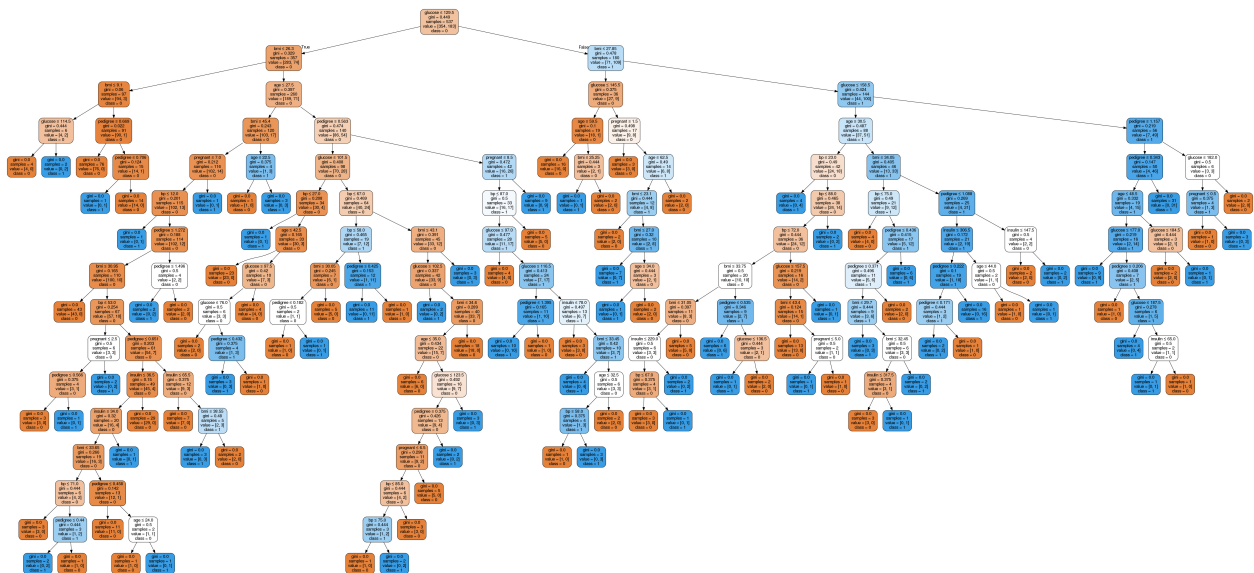
```
pip install graphviz
```

```
pip install pydotplus
```

export_graphviz function converts decision tree classifier into dot file and pydotplus convert this dot file to png or displayable form on Jupyter.

```
from sklearn.tree import export_graphviz
from sklearn.externals.six import StringIO
from IPython.display import Image
import pydotplus

dot_data = StringIO()
export_graphviz(clf, out_file=dot_data,
               filled=True, rounded=True,
               special_characters=True, feature_names = feature_cols, class_names=['0',
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
graph.write_png('diabetes.png')
Image(graph.create_png())
```



[Want to leave a comment?](#)

In the decision tree chart, each internal node has a decision rule that splits the data. Gini referred as Gini ratio, which measures the impurity of the node. You can say a node is pure when all of its records belong to the same class, such nodes known as the leaf node.

Here, the resultant tree is unpruned. This unpruned tree is unexplainable and not easy to understand. In the next section, let's optimize it by pruning.

Optimizing Decision Tree Performance

- **criterion : optional (default="gini") or Choose attribute selection measure:** This parameter allows us to use the different-different attribute selection measure. Supported criteria are "gini" for the Gini index and "entropy" for the information gain.
- **splitter : string, optional (default="best") or Split Strategy:** This parameter allows us to choose the split strategy. Supported strategies are "best" to choose the best split and "random" to choose the best random split.
- **max_depth : int or None, optional (default=None) or Maximum Depth of a Tree:** The maximum depth of the tree. If None, then nodes are expanded until all the leaves contain less than min_samples_split samples. The higher value of maximum depth causes overfitting, and a lower value causes underfitting ([Source](#)).

In Scikit-learn, optimization of decision tree classifier performed by only pre-pruning. Maximum depth of the tree can be used as a control variable for pre-pruning. In the following the example, you can plot a decision tree on the same data with max_depth=3. Other than pre-pruning parameters, You can also try other attribute selection measure such as entropy.

```
# Create Decision Tree classifier object
clf = DecisionTreeClassifier(criterion="entropy", max_depth=3)
```

[Want to leave a comment?](#)

```
# Train Decision Tree Classifier
clf = clf.fit(X_train,y_train)

#Predict the response for test dataset
y_pred = clf.predict(X_test)

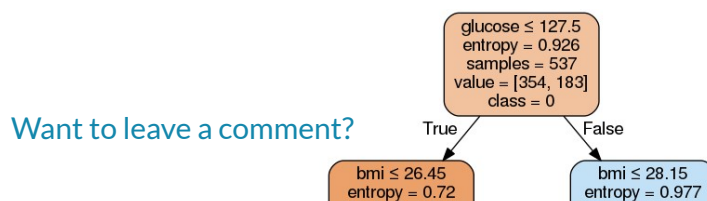
# Model Accuracy, how often is the classifier correct?
print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
```

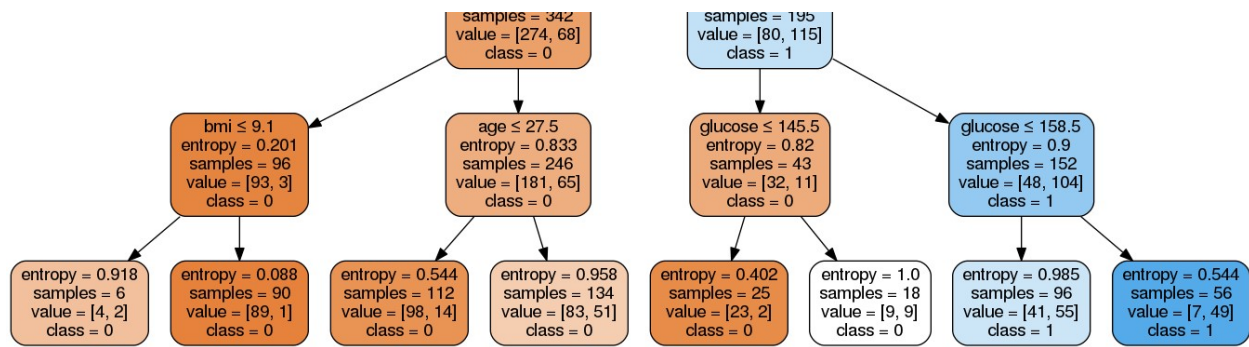
Accuracy: 0.7705627705627706

Well, the classification rate increased to 77.05%, which is better accuracy than the previous model.

Visualizing Decision Trees

```
from sklearn.externals.six import StringIO
from IPython.display import Image
from sklearn.tree import export_graphviz
import pydotplus
dot_data = StringIO()
export_graphviz(clf, out_file=dot_data,
                filled=True, rounded=True,
                special_characters=True, feature_names = feature_cols,class_names=['0', '1'])
graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
graph.write_png('diabetes.png')
Image(graph.create_png())
```





This pruned model is less complex, explainable, and easy to understand than the previous decision tree model plot.

Pros

- Decision trees are easy to interpret and visualize.
- It can easily capture Non-linear patterns.
- It requires fewer data preprocessing from the user, for example, there is no need to normalize columns.
- It can be used for feature engineering such as predicting missing values, suitable for variable selection.
- The decision tree has no assumptions about distribution because of the non-parametric nature of the algorithm. ([Source](#))

Cons

- Sensitive to noisy data. It can overfit noisy data.
- The small variation(or variance) in data can result in the different decision tree. This can be reduced by bagging and boosting algorithms.
- Decision trees are biased with imbalance dataset, so it is recommended that balance [W/out the dataset before?](#) the dataset before creating the decision tree.

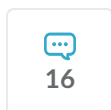
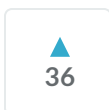
Conclusion

Congratulations, you have made it to the end of this tutorial!

In this tutorial, you covered a lot of details about Decision Tree; It's working, attribute selection measures such as Information Gain, Gain Ratio, and Gini Index, decision tree model building, visualization and evaluation on diabetes dataset using Python Scikit-learn package. Also, discussed its pros, cons, and optimizing Decision Tree performance using parameter tuning.

Hopefully, you can now utilize the Decision tree algorithm to analyze your own datasets.

If you want to learn more about Machine Learning in Python, take DataCamp's [Machine Learning with Tree-Based Models in Python](#) course.



COMMENTS

Dendani Bilal

15/02/2019 10:49 PM

Thank you Avinash.

It was very helpful however I got the following error after I run the code until `y_pred = clf.predict(X_test)`. The following is the error :

ValueError: could not convert string to float: 'DiabetesPedigreeFunction'

▲ 10 ↩ REPLY

[Want to leave a comment?](#)

Avinash Navani

18/03/2019 06:08 AM

▲ 5 ↩ REPLY

Zahid Parvez

18/03/2019 04:35 PM

You need to factorize your categorical data, the model cant deal with non-numeric data.

Have a look at:

- <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.LabelEncoder.html>
- <https://stackoverflow.com/questions/38108832/passing-categorical-data-to-sklearn-decision-tree>



▲ 2 ↩ REPLY

blackstreet23 blackstreet23

07/04/2019 02:30 PM

That still does not solve the problem. The tutorial is buggy ...

▲ 1

Virendra Mishra

11/04/2019 02:02 AM

i got the same error and i used the below script before train the model.

```
IE = preprocessing.LabelEncoder()
```

```
data = yourdatasetname.apply(IE .fit_transform)
```

▲ 1 ↩ REPLY

Colin Contreary

19/04/2019 06:06 PM

If you downloaded the dataset from Kaggle, it comes with a header row. Try replacing it as follows:

```
<code>
```

[Want to leave a comment?](#)

```
col_names = ['pregnant', 'glucose', 'bp', 'skin', 'insulin', 'bmi', 'pedigree', 'age', 'label']
```

```
pima = pd.read_csv("diabetes.csv")
```

```
pima.columns = col_names
```

```
</code>
```

That worked for me!

▲ 3 ↩ [REPLY](#)

Ami Joy

24/04/2019 08:35 AM

Thanks it worked for me also :)

▲ 1

Avinash Navlani

20/04/2019 08:03 PM

[Colin Contreary](#) is right. May be your dataset has columns in it. if you see the data none of the column is String. So need to apply label encoding.

▲ 1 ↩ [REPLY](#)

sahala salim

14/06/2019 12:50 AM

```
ValueError: could not convert string to float: 'Insulin'
```

▲ 1 ↩ [REPLY](#)

kiran kumari

08/04/2019 01:02 PM

Thanks a lot Sir! But I got an error :

```
InvocationException: GraphViz's executables not found
```

Can you please help me out?

[Want to leave a comment?](#)

← **REPLY**

Virendra Mishra

11/04/2019 02:06 AM

you need to install Graphviz, instal pydotplus and setup environment variable. after doing these three step you will see the tree.let me know if you still see error

▲ 2 ← **REPLY**

Colin Contreary

19/04/2019 05:48 PM

I had the same problem of "GraphViz's executables not found." I'm using Windows 10. What worked for me was going into the Environment Variables and adding the graphviz folder to my PATH settings.

You can search "Environment variables" in the search bar of your desktop, or find them manually under Control Panel > System > Advanced System Settings.

Click into Environment Variables and then click Path under the User Variables box.

On my computer, the graphviz folder is stored in
C:\Users\Colin\Anaconda3\Library\bin\graphviz

Find where the graphviz folder is stored on yours, copy that location, then make a new entry in your Path settings with that location.

Then make sure you close down your Jupyter notebook and restart it. Hopefully, the graphviz package will now work!

If not, you can try reading this StackOverflow post in case something else works:

<https://stackoverflow.com/questions/18438997/why-is-pydot-unable-to-find-graphvizs-executables-in-windows-8>

▲ 2 ← **REPLY**

Avinash Navlani

20/04/2019 07:56 PM

Graphviz installation is tricky thing. Initially I also got problem in that. Follow the **Colin Contreary's** Solution.

▲ 1

[Want to leave a comment?](#)

ANURAG NAGAR

I am getting the following error

could not convert string to float:

My data is of string type(catagorical). How could I move ahead?

▲ 1 ↩ REPLY

David Galli

07/05/2019 02:49 AM

State machines are the basics of computer science. It's kind of funny to suggest this as an extra reading material... on the other hand, it's also funny to think that they are a good way to deal with decision trees. Since, if you ever took an [uk essay writing](#) automata class, you'd know that state machines are limited in that they don't have any "memory", which means that you can only describe regular languages using this formalism.

▲ 1 ↩ REPLY

▲ 3

✓ 3

Kevin Ochieng

30/05/2019 10:45 PM

DataCamp kindly try to be clear and procedural like other platforms like towardsdatascience and medium.com ,every time i follow your tutorials errors become my friends and you get other steps skipped without informing the learnings

▲ 1 ↩ REPLY

[Want to leave a comment?](#)

 [Subscribe to RSS](#)



[About](#) [Terms](#) [Privacy](#)

[Want to leave a comment?](#)