

Assignment-15-Random Forest (Company Data)

Random Forest

Assignment

About the data:Let's consider a company dataset with around 10 variables and 400 records.

Sales -- Unit sales (in thousands) at each location Competitor Price -- Price charged by competitor at each location Income -- Community income level (in thousands of dollars) Advertising -- Local advertising budget for company at each location (in thousands of dollars) Population -- Population size in region (in thousands) Price -- Price company charges for car seats at each site Shelf Location at stores -- A factor with levels Bad, Good and Medium indicating the quality of the shelving location for the car seats at each site Age -- Average age of the local population Education -- Education level at each location Urban -- A factor with levels No and Yes to indicate whether the store is in an urban or rural location US -- A factor with levels No and Yes to indicate whether the store is in the US or not The company dataset looks like this:

Problem Statement:

A cloth manufacturing company is interested to know about the segment or attributes causes high sale. Approach - A Random Forest can be built with target variable Sales (we will first convert it in categorical variable) & all other variable will be independent in the analysis.

Import Libraries

In [131]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 %matplotlib inline #for encoding
6 from sklearn.preprocessing import LabelEncoder # for train test splitting,
7 from sklearn.model_selection import train_test_split # for decision tree object
8 from sklearn.tree import DecisionTreeClassifier # for checking testing results
9 from sklearn.metrics import classification_report, confusion_matrix # for visualizing tree
10 from sklearn.tree import plot_tree
11 import warnings
12 warnings.filterwarnings('ignore')
```

UsageError: unrecognized arguments: #for encoding

Import Data

In [132]:

```
1 # Pandas is used for data manipulation
2 import pandas as pd
3 # Read in data and display first 5 rows
4 features = pd.read_csv('Company_Data.csv')
5 features.head(5)
```

Out[132]:

	Sales	CompPrice	Income	Advertising	Population	Price	ShelveLoc	Age	Education	Urban	US
0	9.50	138	73	11	276	120	Bad	42	17	Yes	Yes
1	11.22	111	48	16	260	83	Good	65	10	Yes	Yes
2	10.06	113	35	10	269	80	Medium	59	12	Yes	Yes
3	7.40	117	100	4	466	97	Medium	55	14	Yes	Yes
4	4.15	141	64	3	340	128	Bad	38	13	Yes	No

```
In [133]: 1 # getting information of dataset
          2 features.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 400 entries, 0 to 399
Data columns (total 11 columns):
 #   Column          Non-Null Count  Dtype  
---  --
 0   Sales           400 non-null   float64
 1   CompPrice       400 non-null   int64  
 2   Income          400 non-null   int64  
 3   Advertising     400 non-null   int64  
 4   Population      400 non-null   int64  
 5   Price           400 non-null   int64  
 6   ShelfLoc        400 non-null   object  
 7   Age             400 non-null   int64  
 8   Education       400 non-null   int64  
 9   Urban           400 non-null   object  
10   US              400 non-null   object  
dtypes: float64(1), int64(7), object(3)
memory usage: 34.5+ KB
```

```
In [134]: 1 print('The Shape of our features is:',features.shape)
```

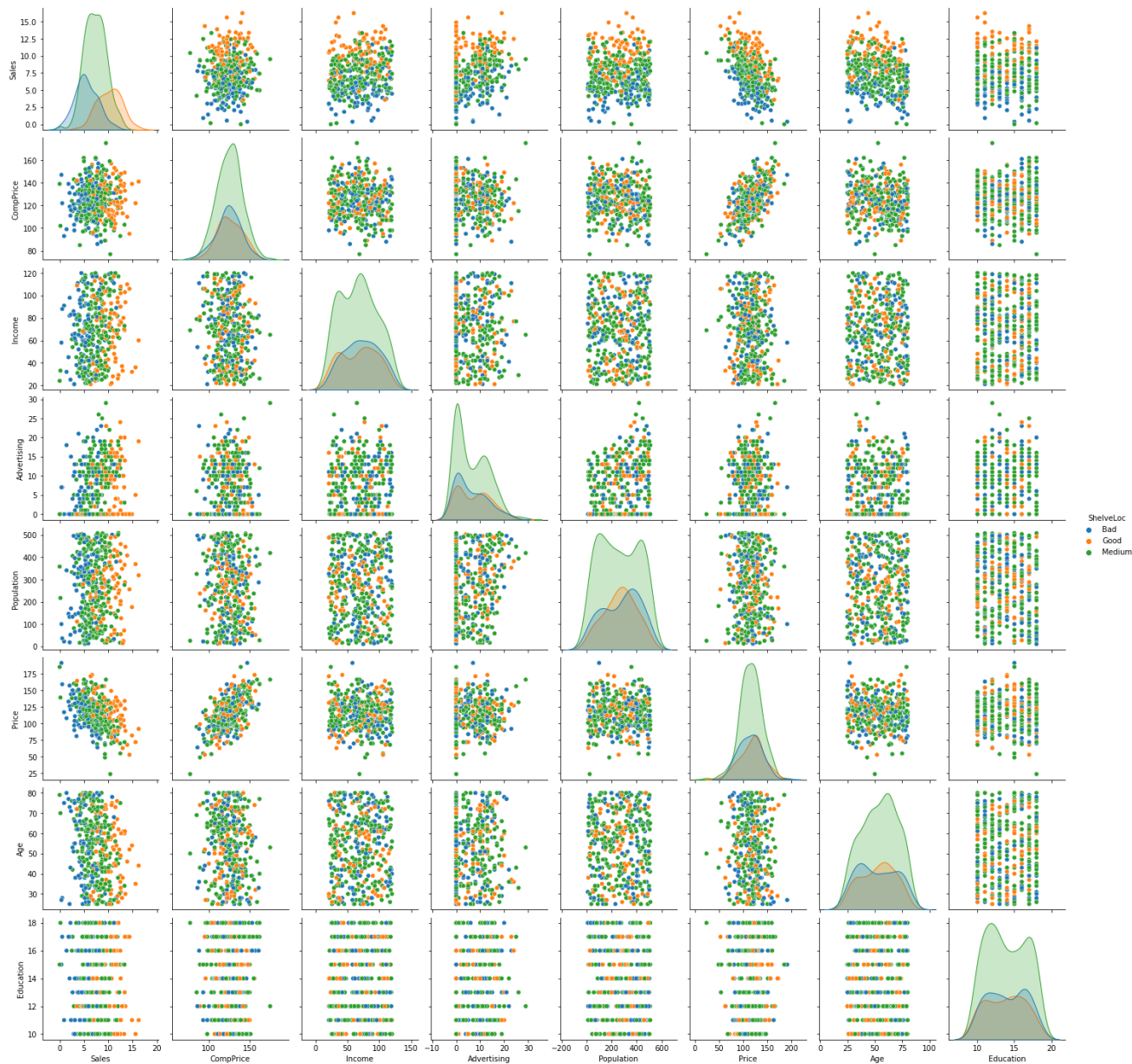
```
The Shape of our features is: (400, 11)
```

```
In [135]: 1 features.isnull().any()
```

```
Out[135]: Sales           False
CompPrice       False
Income          False
Advertising     False
Population      False
Price           False
ShelfLoc        False
Age             False
Education       False
Urban           False
US              False
dtype: bool
```

```
In [136]: 1 # Lets plot pair plot to visualise the attributes all at once.
          2 sns.pairplot(data=features, hue = 'ShelveLoc')
```

Out[136]: <seaborn.axisgrid.PairGrid at 0x227371a58b0>



```
In [137]: 1 # Creating dummy variables dropping first dummy variable
          2 df=pd.get_dummies(features,columns=['Urban','US'], drop_first=True)
```

```
In [138]: 1 print(df.head())
```

	Sales	CompPrice	Income	Advertising	Population	Price	ShelveLoc	Age	\
0	9.50	138	73	11	276	120	Bad	42	
1	11.22	111	48	16	260	83	Good	65	
2	10.06	113	35	10	269	80	Medium	59	
3	7.40	117	100	4	466	97	Medium	55	
4	4.15	141	64	3	340	128	Bad	38	

	Education	Urban_Yes	US_Yes
0	17	1	1
1	10	1	1
2	12	1	1
3	14	1	1
4	13	1	0

```
In [139]: 1 from sklearn.metrics import f1_score
          2 from sklearn.model_selection import train_test_split
```

```
In [140]: 1 df['ShelveLoc']=df['ShelveLoc'].map({'Good':1,'Medium':2,'Bad':3})
```

```
In [141]: 1 print(df.head())
```

	Sales	CompPrice	Income	Advertising	Population	Price	ShelveLoc	Age	\
0	9.50	138	73	11	276	120	3	42	
1	11.22	111	48	16	260	83	1	65	
2	10.06	113	35	10	269	80	2	59	
3	7.40	117	100	4	466	97	2	55	
4	4.15	141	64	3	340	128	3	38	

	Education	Urban_Yes	US_Yes
0	17	1	1
1	10	1	1
2	12	1	1
3	14	1	1
4	13	1	0

```
In [142]: 1 x=df.iloc[:,0:6]
          2 y=df['ShelveLoc']
          3 x
```

```
Out[142]:
```

	Sales	CompPrice	Income	Advertising	Population	Price
0	9.50	138	73	11	276	120
1	11.22	111	48	16	260	83
2	10.06	113	35	10	269	80
3	7.40	117	100	4	466	97
4	4.15	141	64	3	340	128
...
395	12.57	138	108	17	203	128
396	6.14	139	23	3	37	120
397	7.41	162	26	12	368	159
398	5.94	100	79	7	284	95
399	9.71	134	37	0	27	120

400 rows × 6 columns

```
In [143]: 1 y
```

```
Out[143]: 0    3
          1    1
          2    2
          3    2
          4    3
          ..
          395    1
          396    2
          397    2
          398    3
          399    1
          Name: ShelveLoc, Length: 400, dtype: int64
```

```
In [144]: 1 df['ShelveLoc'].unique()
```

```
Out[144]: array([3, 1, 2], dtype=int64)
```

```
In [145]: 1 df.ShelveLoc.value_counts()
```

```
Out[145]: 2    219
          3    96
          1    85
          Name: ShelveLoc, dtype: int64
```

```
In [146]: 1 colnames = list(df.columns)
          2 colnames
```

```
Out[146]: ['Sales',
           'CompPrice',
           'Income',
           'Advertising',
           'Population',
           'Price',
           'ShelveLoc',
           'Age',
           'Education',
           'Urban_Yes',
           'US_Yes']
```

```
In [147]: 1 # Descriptive statistics for each column
          2 df.describe()
```

```
Out[147]:
```

	Sales	CompPrice	Income	Advertising	Population	Price	ShelveLoc	Age	Education	Urban_Yes	US_Yes
count	400.000000	400.000000	400.000000	400.000000	400.000000	400.000000	400.000000	400.000000	400.000000	400.000000	400.000000
mean	7.496325	124.975000	68.657500	6.635000	264.840000	115.795000	2.027500	53.322500	13.900000	0.705000	0.645000
std	2.824115	15.334512	27.986037	6.650364	147.376436	23.676664	0.672961	16.200297	2.620528	0.456614	0.479113
min	0.000000	77.000000	21.000000	0.000000	10.000000	24.000000	1.000000	25.000000	10.000000	0.000000	0.000000
25%	5.390000	115.000000	42.750000	0.000000	139.000000	100.000000	2.000000	39.750000	12.000000	0.000000	0.000000
50%	7.490000	125.000000	69.000000	5.000000	272.000000	117.000000	2.000000	54.500000	14.000000	1.000000	1.000000
75%	9.320000	135.000000	91.000000	12.000000	398.500000	131.000000	2.000000	66.000000	16.000000	1.000000	1.000000
max	16.270000	175.000000	120.000000	29.000000	509.000000	191.000000	3.000000	80.000000	18.000000	1.000000	1.000000

```
In [148]: 1 df.head()
```

```
Out[148]:
```

	Sales	CompPrice	Income	Advertising	Population	Price	ShelveLoc	Age	Education	Urban_Yes	US_Yes
0	9.50	138	73	11	276	120	3	42	17	1	1
1	11.22	111	48	16	260	83	1	65	10	1	1
2	10.06	113	35	10	269	80	2	59	12	1	1
3	7.40	117	100	4	466	97	2	55	14	1	1
4	4.15	141	64	3	340	128	3	38	13	1	0

```
In [149]: 1 # Labels are the values we want to predict
          2 labels = np.array(df['Income'])
          3 # Remove the Labels from the features
          4 # axis 1 refers to the columns
          5 features=df.drop('Income', axis = 1)
          6 # Saving feature names for later use
          7 features_list = list(df.columns)
          8 # Convert to numpy array
          9 features = np.array(df)
```

```
In [150]: 1 # Using Skicit-Learn to split data into training and testing sets
          2 from sklearn.model_selection import train_test_split
          3 # Split the data into training and testing sets
          4 train_features, test_features, train_labels, test_labels = train_test_split(features, labels, test_size = 0.25, random_state
```

```
In [151]: 1 print('Training Features Shape:', train_features.shape)
          2 print('Training Labels Shape:', train_labels.shape)
          3 print('Testing Features Shape:', test_features.shape)
          4 print('Testing Labels Shape:', test_labels.shape)
```

```
Training Features Shape: (300, 11)
Training Labels Shape: (300,)
Testing Features Shape: (100, 11)
Testing Labels Shape: (100,)
```

Establish Baseline

```
In [152]: 1 # The baseline predictions are the historical averages
2 baseline_preds = test_features[:, features_list.index('Sales')]
3 # Baseline errors, and display average baseline error
4 baseline_errors = abs(baseline_preds - test_labels)
5 print('Average baseline error: ', round(np.mean(baseline_errors), 2))
```

Average baseline error: 65.26

```
In [153]: 1 # Import the model we are using
2 from sklearn.ensemble import RandomForestRegressor
3 # Instantiate model with 1000 decision trees
4 rf = RandomForestRegressor(n_estimators = 1000, random_state = 42)
5 # Train the model on training data
6 rf.fit(train_features, train_labels);
```

```
In [154]: 1 # Use the forest's predict method on the test data
2 predictions = rf.predict(test_features)
3 # Calculate the absolute errors
4 errors = abs(predictions - test_labels)
5 # print out the mean absolute error (mae)
6 print('Meam Absolute Error:', round(np.mean(errors), 2), 'degrees.')
```

Meam Absolute Error: 0.27 degrees.

Determine Performance Metrics

```
In [155]: 1 # Calculate mean absolute percentage error (MAPE)
2 mape = 100 * (errors / test_labels)
3 # Calculate and display accuracy
4 accuracy = 100 - np.mean(mape)
5 print('Accuracy:', round(accuracy, 2), '%.')
```

Accuracy: 99.58 %.

```
In [156]: 1 !pip install pydot
```

Requirement already satisfied: pydot in c:\users\admin\anaconda3\lib\site-packages (1.4.2)

Requirement already satisfied: pyparsing>=2.1.4 in c:\users\admin\anaconda3\lib\site-packages (from pydot) (3.0.4)

```
In [157]: 1 !pip install graphviz
```

Requirement already satisfied: graphviz in c:\users\admin\anaconda3\lib\site-packages (0.20.1)

```
In [171]: 1 graph.write_png
```

```
Out[171]: <function pydot.Dot.__init__.<locals>.new_method(path, f='png', prog='dot', encoding=None)>
```

```
In [169]: 1 # Import tools needed for visualization
2 from sklearn.tree import export_graphviz
3 import pydot
4 # Pull out one tree from the forest
5 tree = rf.estimators_[5]
6 # Import tools needed for visualization
7 from sklearn.tree import export_graphviz
8 import pydot
9 # Pull out one tree from the forest
10 tree = rf.estimators_[5]
11 # Export the image to a dot file
12 export_graphviz(tree, out_file = 'tree.dot', feature_names = features_list, rounded = True, precision = 1)
13 # Use dot file to create a graph
14 (graph, ) = pydot.graph_from_dot_file('tree.dot')
15 # Write graph to a png file
16 graph.write_png('tree.png')
```

```
-----
FileNotFoundError                                Traceback (most recent call last)
~\anaconda3\lib\site-packages\pydot.py in create(self, prog, format, encoding)
    1922         try:
-> 1923             stdout_data, stderr_data, process = call_graphviz(
    1924                 program=prog,

~\anaconda3\lib\site-packages\pydot.py in call_graphviz(program, arguments, working_dir, **kwargs)
    131
-> 132     process = subprocess.Popen(
    133         program_with_args,

~\anaconda3\lib\subprocess.py in __init__(self, args, bufsize, executable, stdin, stdout, stderr, preexec_fn, close_fds, shell,
cwd, env, universal_newlines, startupinfo, creationflags, restore_signals, start_new_session, pass_fds, user, group, extra_grou
ps, encoding, errors, text, umask)
    950
-> 951         self._execute_child(args, executable, preexec_fn, close_fds,
    952                             pass_fds, cwd, env,

~\anaconda3\lib\subprocess.py in _execute_child(self, args, executable, preexec_fn, close_fds, pass_fds, cwd, env, startupinfo,
creationflags, shell, p2cread, p2cwrite, c2pread, c2pwrite, errread, errwrite, unused_restore_signals, unused_gid, unused_gids,
unused_uid, unused_umask, unused_start_new_session)
    1419         try:
-> 1420             hp, ht, pid, tid = _winapi.CreateProcess(executable, args,
    1421                                                         # no special security

FileNotFoundError: [WinError 2] The system cannot find the file specified
```

During handling of the above exception, another exception occurred:

```
FileNotFoundError                                Traceback (most recent call last)
~\AppData\Local\Temp\ipykernel_6664\2350536375.py in <module>
     14 (graph, ) = pydot.graph_from_dot_file('tree.dot')
     15 # Write graph to a png file
--> 16 graph.write_png('tree.png')

~\anaconda3\lib\site-packages\pydot.py in new_method(path, f, prog, encoding)
    1741         encoding=None):
    1742             """Refer to docstring of method `write`."""
-> 1743         self.write(
    1744             path, format=f, prog=prog,
    1745             encoding=encoding)

~\anaconda3\lib\site-packages\pydot.py in write(self, path, prog, format, encoding)
    1826         f.write(s)
    1827     else:
-> 1828         s = self.create(prog, format, encoding=encoding)
    1829         with io.open(path, mode='wb') as f:
    1830             f.write(s)

~\anaconda3\lib\site-packages\pydot.py in create(self, prog, format, encoding)
    1931         args[1] = "{prog}" not found in path.'.format(
    1932             prog=prog)
-> 1933         raise OSError(*args)
    1934     else:
    1935         raise

FileNotFoundError: [WinError 2] "dot" not found in path.
```

```

In [166]: 1 # Import tools needed for visualization
2 from sklearn.tree import export_graphviz
3 import pydot
4 import graphviz
5 # Pull out one tree from the forest
6 tree = rf.estimators_[5]
7 # Import tools needed for visualization
8 from sklearn.tree import export_graphviz
9 import pydot
10 # Pull out one tree from the forest
11 tree = rf.estimators_[5]
12 # Export the image to a dot file
13 export_graphviz(tree, out_file = 'tree.dot', feature_names = features_list, rounded = True, precision = 1)
14 # Use dot file to create a graph
15 (graph, ) = pydot.graph_from_dot_file('tree.dot')
16 # Write graph to a png file
17 graph.write_png('tree.png')

```

```

FileNotFoundError                                Traceback (most recent call last)
~\anaconda3\lib\site-packages\pydot.py in create(self, prog, format, encoding)
    1922         try:
-> 1923             stdout_data, stderr_data, process = call_graphviz(
    1924                 program=prog,

~\anaconda3\lib\site-packages\pydot.py in call_graphviz(program, arguments, working_dir, **kwargs)
    131
-> 132     process = subprocess.Popen(
    133         program_with_args,

~\anaconda3\lib\subprocess.py in __init__(self, args, bufsize, executable, stdin, stdout, stderr, preexec_fn, close_fds, shell,
cwd, env, universal_newlines, startupinfo, creationflags, restore_signals, start_new_session, pass_fds, user, group, extra_grou
ps, encoding, errors, text, umask)
    950
-> 951         self._execute_child(args, executable, preexec_fn, close_fds,
    952                             pass_fds, cwd, env,

~\anaconda3\lib\subprocess.py in _execute_child(self, args, executable, preexec_fn, close_fds, pass_fds, cwd, env, startupinfo,
creationflags, shell, p2cread, p2cwrite, c2pread, c2pwrite, errread, errwrite, unused_restore_signals, unused_gid, unused_gids,
unused_uid, unused_umask, unused_start_new_session)
    1419         try:
-> 1420             hp, ht, pid, tid = _winapi.CreateProcess(executable, args,
    1421                                                         # no special security

FileNotFoundError: [WinError 2] The system cannot find the file specified

During handling of the above exception, another exception occurred:

FileNotFoundError                                Traceback (most recent call last)
~\AppData\Local\Temp\ipykernel_6664\2479579253.py in <module>
     15 (graph, ) = pydot.graph_from_dot_file('tree.dot')
     16 # Write graph to a png file
----> 17 graph.write_png('tree.png')

~\anaconda3\lib\site-packages\pydot.py in new_method(path, f, prog, encoding)
    1741         encoding=None):
    1742         """Refer to docstring of method `write`."""
-> 1743         self.write(
    1744             path, format=f, prog=prog,
    1745             encoding=encoding)

~\anaconda3\lib\site-packages\pydot.py in write(self, path, prog, format, encoding)
    1826         f.write(s)
    1827     else:
-> 1828         s = self.create(prog, format, encoding=encoding)
    1829         with io.open(path, mode='wb') as f:
    1830             f.write(s)

~\anaconda3\lib\site-packages\pydot.py in create(self, prog, format, encoding)
    1931         args[1] = "{prog}" not found in path.'.format(
    1932             prog=prog)
-> 1933         raise OSError(*args)
    1934     else:
    1935         raise

FileNotFoundError: [WinError 2] ".dot" not found in path.

```



```
In [159]: 1 # Limit depth of tree to 3 levels
2 rf_small = RandomForestRegressor(n_estimators=10, max_depth = 3)
3 rf_small.fit(train_features, train_labels)
4 # Extract the small tree
5 tree_small = rf_small.estimators_[5]
6 # Save the tree as a png image
7 export_graphviz(tree_small, out_file = 'small_tree.dot', feature_names = features_list, rounded = True, precision = 1)
8 (graph, ) = pydot.graph_from_dot_file('small_tree.dot')
9 graph.write_png('small_tree.png');
```

```
-----
FileNotFoundError Traceback (most recent call last)
~\anaconda3\lib\site-packages\pydot.py in create(self, prog, format, encoding)
    1922         try:
-> 1923             stdout_data, stderr_data, process = call_graphviz(
    1924                 program=prog,

~\anaconda3\lib\site-packages\pydot.py in call_graphviz(program, arguments, working_dir, **kwargs)
    131
-> 132     process = subprocess.Popen(
    133         program_with_args,

~\anaconda3\lib\subprocess.py in __init__(self, args, bufsize, executable, stdin, stdout, stderr, preexec_fn, close_fds, shell,
cwd, env, universal_newlines, startupinfo, creationflags, restore_signals, start_new_session, pass_fds, user, group, extra_grou
ps, encoding, errors, text, umask)
    950
-> 951         self._execute_child(args, executable, preexec_fn, close_fds,
    952                             pass_fds, cwd, env,

~\anaconda3\lib\subprocess.py in _execute_child(self, args, executable, preexec_fn, close_fds, pass_fds, cwd, env, startupinfo,
creationflags, shell, p2cread, p2cwrite, c2pread, c2pwrite, errread, errwrite, unused_restore_signals, unused_gid, unused_gids,
unused_uid, unused_umask, unused_start_new_session)
    1419         try:
-> 1420             hp, ht, pid, tid = _winapi.CreateProcess(executable, args,
    1421                                                         # no special security
```

FileNotFoundError: [WinError 2] The system cannot find the file specified

During handling of the above exception, another exception occurred:

```
FileNotFoundError Traceback (most recent call last)
~\AppData\Local\Temp\ipykernel_6664\1441242784.py in <module>
      7 export_graphviz(tree_small, out_file = 'small_tree.dot', feature_names = features_list, rounded = True, precision = 1)
      8 (graph, ) = pydot.graph_from_dot_file('small_tree.dot')
----> 9 graph.write_png('small_tree.png');
```

```
~\anaconda3\lib\site-packages\pydot.py in new_method(path, f, prog, encoding)
    1741         encoding=None):
    1742             """Refer to docstring of method `write`."""
-> 1743         self.write(
    1744             path, format=f, prog=prog,
    1745             encoding=encoding)

~\anaconda3\lib\site-packages\pydot.py in write(self, path, prog, format, encoding)
    1826         f.write(s)
    1827     else:
-> 1828         s = self.create(prog, format, encoding=encoding)
    1829         with io.open(path, mode='wb') as f:
    1830             f.write(s)

~\anaconda3\lib\site-packages\pydot.py in create(self, prog, format, encoding)
    1931         args[1] = "{prog}" not found in path.'.format(
    1932             prog=prog)
-> 1933         raise OSError(*args)
    1934     else:
    1935         raise
```

FileNotFoundError: [WinError 2] "dot" not found in path.

```
In [161]: 1 # Get numerical feature importances
2 importances = list(rf.feature_importances_)
3 # List of tuples with variable and importance
4 feature_importances = [(feature, round(importance, 2)) for feature, importance in zip(features_list, importances)]
5 # Sort the feature importances by most important first
6 feature_importances = sorted(feature_importances, key = lambda x: x[1], reverse = True)
7 # Print out the feature and importances
8 [print('Variable: {} Importance: {}'.format(*pair)) for pair in feature_importances];
```

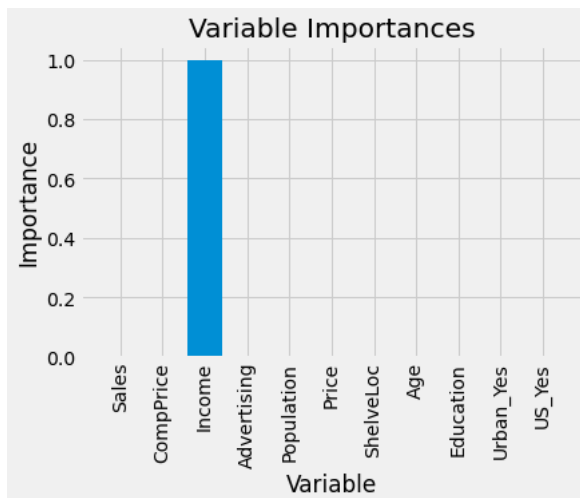
```
Variable: Income           Importance: 1.0
Variable: Sales            Importance: 0.0
Variable: CompPrice        Importance: 0.0
Variable: Advertising      Importance: 0.0
Variable: Population       Importance: 0.0
Variable: Price            Importance: 0.0
Variable: Shelveloc        Importance: 0.0
Variable: Age              Importance: 0.0
Variable: Education        Importance: 0.0
Variable: Urban_Yes        Importance: 0.0
Variable: US_Yes           Importance: 0.0
```

```
In [163]: 1 # New random forest with only the two most important variables
2 rf_most_important = RandomForestRegressor(n_estimators=1000, random_state=42)
3 # Extract the two most important features
4 important_indices = [features_list.index('Sales'), features_list.index('Income')]
5 train_important = train_features[:, important_indices]
6 test_important = test_features[:, important_indices]
7 # Train the random forest
8 rf_most_important.fit(train_important, train_labels)
9 # Make predictions and determine the error
10 predictions = rf_most_important.predict(test_important)
11 errors = abs(predictions - test_labels)
12 # Display the performance metrics
13 print('Mean Absolute Error:', round(np.mean(errors), 2), 'degrees.')
14 mape = np.mean(100 * (errors / test_labels))
15 accuracy = 100 - mape
16 print('Accuracy:', round(accuracy, 2), '%.')
```

```
Mean Absolute Error: 0.16 degrees.
Accuracy: 99.75 %.
```

```
In [165]: 1 # Import matplotlib for plotting and use magic command for Jupyter Notebooks
2 import matplotlib.pyplot as plt
3 %matplotlib inline
4 # Set the style
5 plt.style.use('fivethirtyeight')
6 # List of x locations for plotting
7 x_values = list(range(len(importances)))
8 # Make a bar chart
9 plt.bar(x_values, importances, orientation = 'vertical')
10 # Tick labels for x axis
11 plt.xticks(x_values, features_list, rotation='vertical')
12 # Axis labels and title
13 plt.ylabel('Importance'); plt.xlabel('Variable'); plt.title('Variable Importances')
```

```
Out[165]: Text(0.5, 1.0, 'Variable Importances')
```



In []:

1

In []:

1