## ▾ Importing Libraries

```
1 import numpy as np
2 import pandas as pd
3
4 # Set random seed to ensure reproducible runs
5 RSEED = 50
```

We create the features X and the labels y. There are only two features, which will allow us to visualize the data and which makes this a very easy problem.

```
1
2 X = np.array([[2, 2],
3               [2, 1],
4               [2, 3],
5               [1, 2],
6               [1, 1],
7               [3, 3]])
8
9 y = np.array([0, 1, 1, 1, 0, 1])
```
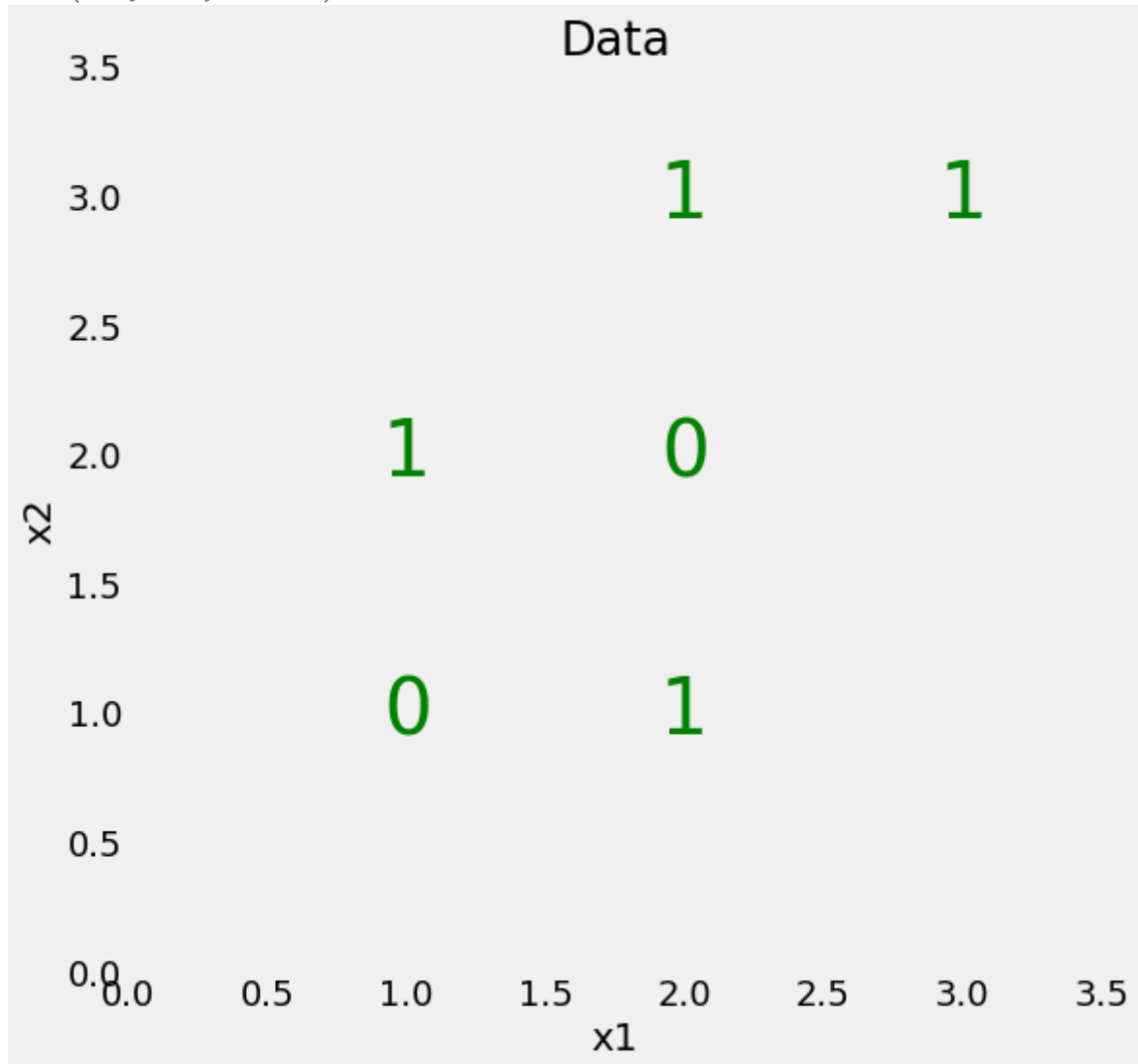
## ▾ Data Visualization

To get a sense of the data, we can graph the data points with the number showing the label.

```
1  import matplotlib.pyplot as plt
2  %matplotlib inline
3  import seaborn as sns
4
5  # Plot formatting
6  plt.style.use('fivethirtyeight')
7  plt.rcParams['font.size'] = 18
8  plt.figure(figsize = (8, 8))
9
10 # Plot each point as the label
11 for x1, x2, label in zip(X[:, 0], X[:, 1], y):
12     plt.text(x1, x2, str(label), fontsize = 40, color = 'g',
13              ha='center', va='center')
14
15 # Plot formatting
16 plt.grid(None);
17 plt.xlim((0, 3.5));
18 plt.ylim((0, 3.5));
19 plt.xlabel('x1', size = 20); plt.ylabel('x2', size = 20); plt.title('Data', size = 24)
```

```
/usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarni
  import pandas.util.testing as tm
Text(0.5, 1.0, 'Data')
```



## Decision Tree

Here we quickly build and train a single decision tree on the data using Scikit-Learn. The tree will
learn how to separate the points, building a flowchart of questions based on the feature values
and the labels. At each stage, the decision tree makes splits by maximizing the reduction in Gini
impurity.

We'll use the default hyperparameters for the decision tree which means it can grow as deep as
necessary in order to completely separate the classes. This will lead to overfitting because the
model memorizes the training data, and in practice, we usually want to limit the depth of the
tree so it can generalize to testing data.

```
1 from sklearn.tree import DecisionTreeClassifier
2
3 # Make a decision tree and train
4 tree = DecisionTreeClassifier(random_state=RSEED)
5 tree.fit(X, y)
```

```
DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                       max_depth=None, max_features=None, max_leaf_nodes=None,
                       min_impurity_decrease=0.0, min_impurity_split=None,
                       min_samples_leaf=1, min_samples_split=2,
                       min_weight_fraction_leaf=0.0, presort='deprecated',
                       random state=50  splitter='best')
```

```
1 print(f'Decision tree has {tree.tree_.node_count} nodes with maximum depth {tree.tree_.
```

```
Decision tree has 9 nodes with maximum depth 3.
```

```
1 print(f'Model Accuracy: {tree.score(X, y)}')
```

```
Model Accuracy: 1.0
```

## Visualize Decision Tree

To get a sense of how the decision tree "thinks", it's helpful to visualize the entire structure. This will show each node in the tree which we can use to make new predictions. Because the tree is relatively small, we can understand the entire image.

## Gini Impurity

The Gini Impurity represents the probability that a randomly selected sample from the node will be incorrectly classified according to the distribution of samples in the node. At the top, there is a 44.4% chance that a randomly selected point would be incorrectly classified. The Gini Impurity is how the decision tree makes splits. It splits the samples based on the value of a feature that reduces the Gini Impurity by the largest amount. If we do the math, the average (weighted by number of samples) Gini Impurity decreases as we move down the tree.

Eventually, the average Gini Impurity goes to 0.0 as we correctly classify each point. However, correctly classifying every single training point is usually not a good indicator because that means the model will not be able to generalize to the testing data! This model correclty classifies every single training point because we did not limit the maximum depth and during training, we give the model the answers as well as the features.

First we export the tree as a dot file making sure to label the features and the classes.

```
1 from sklearn.tree import export_graphviz
2
3 # Export as dot
4 export_graphviz(tree, 'tree.dot', rounded = True,
5                 feature_names = ['x1', 'x2'],
6                 class_names = ['0', '1'], filled = True)
```

```
1 from subprocess import call
2 # Convert to png
3 call(['dot', '-Tpng', 'tree.dot', '-o', 'tree.png', '-Gdpi=400']);
```

```
1
2 from IPython.display import Image
3 Image('tree.png')
```
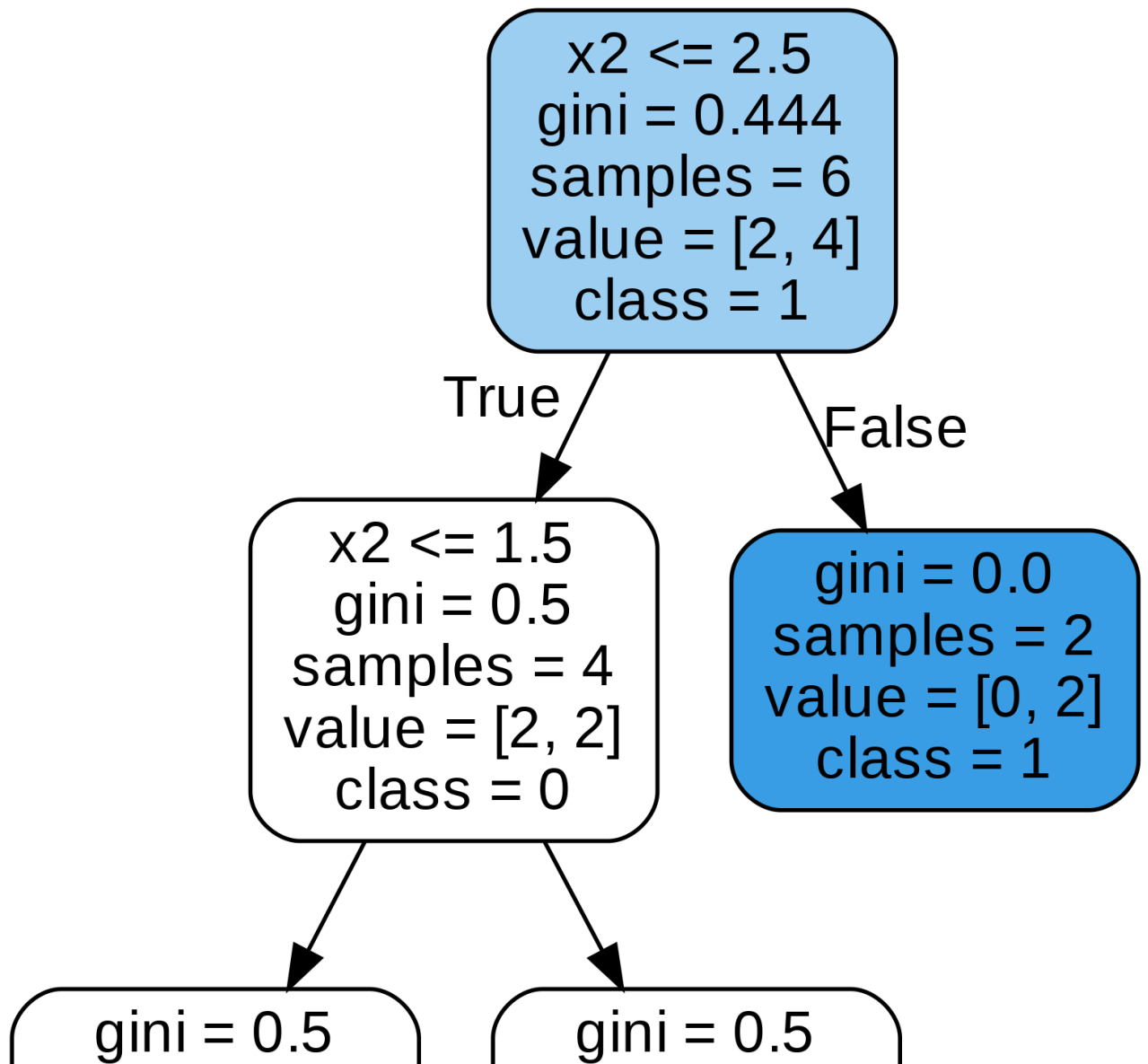
x2 <= 2.5
gini = 0.444
samples = 6

```
1 # Limit maximum depth and train
2 short_tree = DecisionTreeClassifier(max_depth = 2, random_state=RSEED)
3 short_tree.fit(X, y)
4
5 print(f'Model Accuracy: {short_tree.score(X, y)}')
```

⊳  Model Accuracy: 0.6666666666666666

samples = 4          value = [0, 2]

```
1
2 # Export as dot
3 export_graphviz(short_tree, 'shorttree.dot', rounded = True,
4                 feature_names = ['x1', 'x2'],
5                 class_names = ['0', '1'], filled = True)
6
7 call(['dot', '-Tpng', 'shorttree.dot', '-o', 'shorttree.png', '-Gdpi=400']);
8 Image('shorttree.png')
```

⊳

```
x2 <= 2.5
gini = 0.444
samples = 6
value = [2, 4]
class = 1
```

True

False

```
x2 <= 1.5
gini = 0.5
samples = 4
value = [2, 2]
class = 0
```

```
gini = 0.0
samples = 2
value = [0, 2]
class = 1
```

```
gini = 0.5
```

```
gini = 0.5
```

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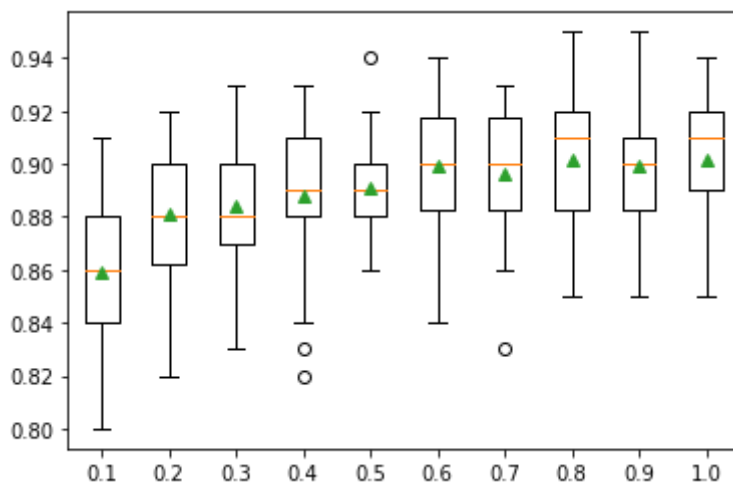
```
24        i = None
25     models[key] = RandomForestClassifier(max_samples=i)
26   return models
27
28 # evaluate a given model using cross-validation
29 def evaluate_model(model, X, y):
30   # define the evaluation procedure
31   cv = RepeatedStratifiedKFold(n_splits=10, n_repeats=3, random_state=1)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

```
>0.1 0.859 (0.028)
>0.2 0.881 (0.025)
>0.3 0.884 (0.021)
>0.4 0.888 (0.028)
>0.5 0.891 (0.019)
>0.6 0.900 (0.023)
>0.7 0.897 (0.024)
>0.8 0.901 (0.024)
>0.9 0.900 (0.026)
>1.0 0.901 (0.026)
```
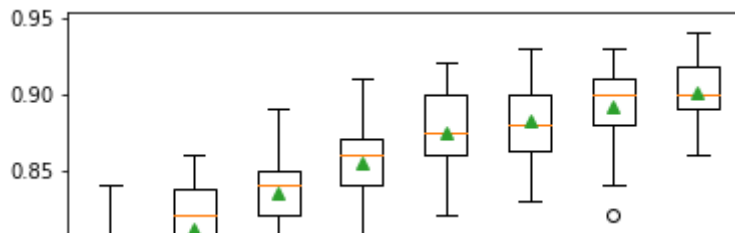


```
1 # explore random forest tree depth effect on performance
2 from numpy import mean
```

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

```
>1 0.766 (0.045)
>2 0.811 (0.039)
>3 0.836 (0.032)
>4 0.854 (0.027)
>5 0.874 (0.028)
>6 0.882 (0.026)
>7 0.891 (0.027)
>None 0.901 (0.022)
```



## Random Forest Classifer using a dataset

```
1 #Start by connecting gdrive into the google colab
2
3 from google.colab import drive
4
5 drive.mount('/content/gdrive')
```

```
Mounted at /content/gdrive
```

```
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('/content/gdrive/My Drive/Colab Notebooks/clustering/temps.csv')
5 features.head(5)
```

|   | year | month | day | week | temp_2 | temp_1 | average | actual | forecast_noaa | forecast_ |
|---|------|-------|-----|------|--------|--------|---------|--------|---------------|-----------|
| 0 | 2016 | 1 | 1 | Fri | 45 | 45 | 45.6 | 45 | 43 | |
| 1 | 2016 | 1 | 2 | Sat | 44 | 45 | 45.7 | 44 | 41 | |
| 2 | 2016 | 1 | 3 | Sun | 45 | 44 | 45.8 | 41 | 43 | |
| 3 | 2016 | 1 | 4 | Mon | 44 | 41 | 45.9 | 40 | 44 | |
| 4 | 2016 | 1 | 5 | Tues | 41 | 40 | 46.0 | 44 | 46 | |

```
1 print('The shape of our features is:', features.shape)
2
```

```
The shape of our features is: (348, 12)
```

```
1 # Descriptive statistics for each column
2 features.describe()
```

```

```

|  | year | month | day | temp_2 | temp_1 | average | actual |
|---|---|---|---|---|---|---|---|
| count | 348.0 | 348.000000 | 348.000000 | 348.000000 | 348.000000 | 348.000000 | 348.000000 |
| mean | 2016.0 | 6.477011 | 15.514368 | 62.652299 | 62.701149 | 59.760632 | 62.543103 |
| std | 0.0 | 3.498380 | 8.772982 | 12.165398 | 12.120542 | 10.527306 | 11.794146 |
| min | 2016.0 | 1.000000 | 1.000000 | 35.000000 | 35.000000 | 45.100000 | 35.000000 |
| 25% | 2016.0 | 3.000000 | 8.000000 | 54.000000 | 54.000000 | 49.975000 | 54.000000 |
| 50% | 2016.0 | 6.000000 | 15.000000 | 62.500000 | 62.500000 | 58.200000 | 62.500000 |
| 75% | 2016.0 | 10.000000 | 22.000000 | 71.000000 | 71.000000 | 69.025000 | 71.000000 |

```
1 # One-hot encode the data using pandas get_dummies
2 features = pd.get_dummies(features)
3 # Display the first 5 rows of the last 12 columns
4 features.iloc[:,5:].head(5)
```

|  | average | actual | forecast_noaa | forecast_acc | forecast_under | friend | week_Fri | w |
|---|---|---|---|---|---|---|---|---|
| 0 | 45.6 | 45 | 43 | 50 | 44 | 29 | 1 | |
| 1 | 45.7 | 44 | 41 | 50 | 44 | 61 | 0 | |
| 2 | 45.8 | 41 | 43 | 46 | 47 | 56 | 0 | |
| 3 | 45.9 | 40 | 44 | 48 | 46 | 53 | 0 | |
| 4 | 46.0 | 44 | 46 | 46 | 46 | 41 | 0 | |

```
1 # Use numpy to convert to arrays
2 import numpy as np
3 # Labels are the values we want to predict
4 labels = np.array(features['actual'])
5 # Remove the labels from the features
6 # axis 1 refers to the columns
7 features= features.drop('actual', axis = 1)
8 # Saving feature names for later use
9 feature_list = list(features.columns)
10 # Convert to numpy array
11 features = np.array(features)
```

```
1 # Using Scikit-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, l
5
```

```
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: (261, 17)
Training Labels Shape: (261,)
Testing Features Shape: (87, 17)
Testing Labels Shape: (87,)
```

```
1 # The baseline predictions are the historical averages
2 baseline_preds = test_features[:, feature_list.index('average')]
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:  5.06
```

```
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);
```

```
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('Mean Absolute Error:', round(np.mean(errors), 2), 'degrees.')
```

```
Mean Absolute Error: 3.87 degrees.
```

```
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), '%.')
6
```
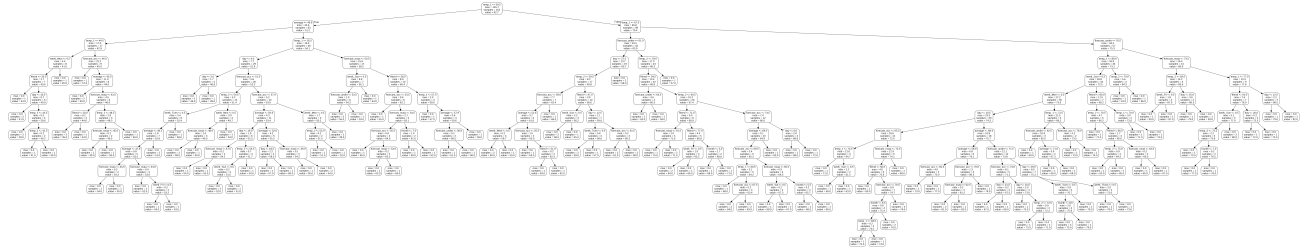
```
Accuracy: 93.93 %.
```

```
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 = 'tree1.dot', feature_names = feature_list, rounded = 1
13 # Use dot file to create a graph
14 (graph, ) = pydot.graph_from_dot_file('tree1.dot')
15 # Write graph to a png file
```

```
16 graph.write_png('tree1.png')
17 from IPython.display import Image
18 Image('tree1.png')
```
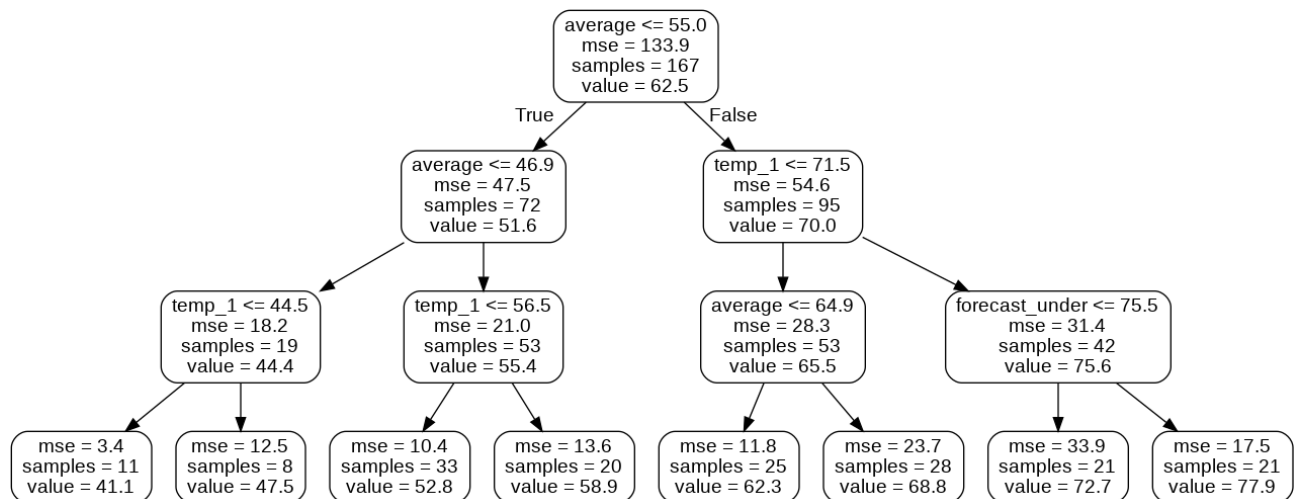


```
1 # Limit depth of tree to 3 levels
2 rf_small = RandomForestRegressor(n_estimators=10, max_depth = 3)
3 rf_small.fit(train_features, train_labels)
4 # Extract the small tree
5 tree_small = rf_small.estimators_[5]
6 # Save the tree as a png image
7 export_graphviz(tree_small, out_file = 'small_tree.dot', feature_names = feature_list,
8 (graph, ) = pydot.graph_from_dot_file('small_tree.dot')
9 graph.write_png('small_tree.png');
10 from IPython.display import Image
11 Image('small_tree.png')
12
13
```



```
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(f
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: {:20} Importance: {}'.format(*pair)) for pair in feature_importances]
```

```
Variable: temp_1               Importance: 0.66
Variable: average              Importance: 0.15
Variable: forecast_noaa        Importance: 0.05
Variable: forecast_acc         Importance: 0.03
Variable: day                  Importance: 0.02
Variable: temp_2               Importance: 0.02
Variable: forecast_under       Importance: 0.02
Variable: friend               Importance: 0.02
Variable: month                Importance: 0.01
Variable: year                 Importance: 0.0
Variable: week_Fri             Importance: 0.0
Variable: week_Mon             Importance: 0.0
Variable: week_Sat             Importance: 0.0
Variable: week_Sun             Importance: 0.0
Variable: week_Thurs           Importance: 0.0
Variable: week_Tues            Importance: 0.0
Variable: week_Wed             Importance: 0.0
```

Double-click (or enter) to edit

```
 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 = [feature_list.index('temp_1'), feature_list.index('average')]
 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: 3.92 degrees.
Accuracy: 93.76 %.
```
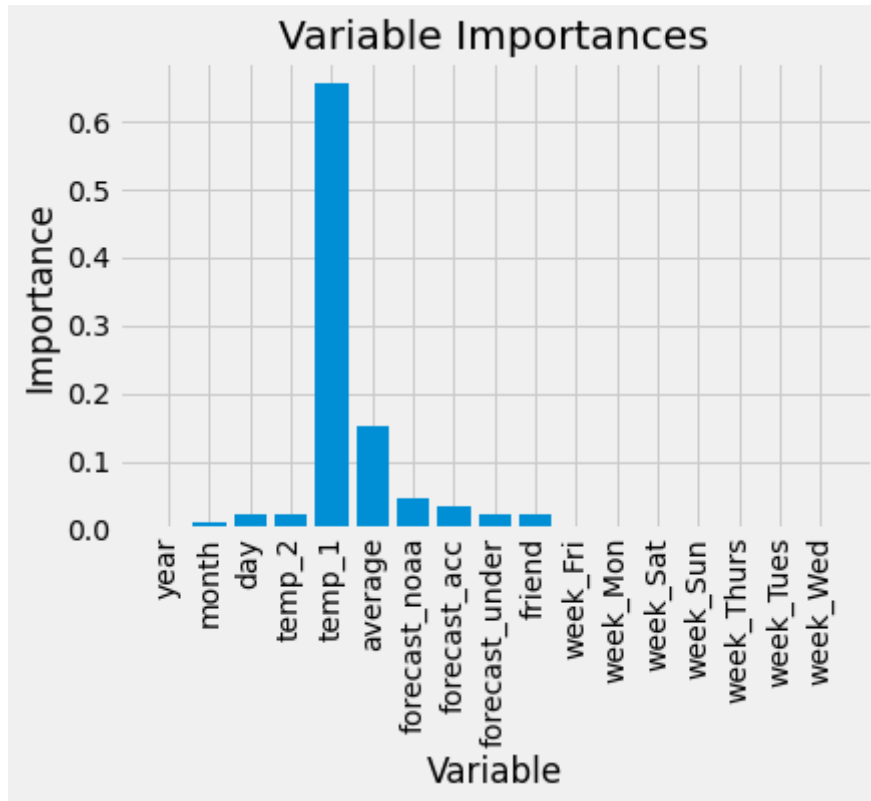
```
 1 # Import matplotlib for plotting and use magic command for Jupyter Notebooks
 2 import matplotlib.pyplot as plt
 3 %matplotlib inline
 4 # Set the style
 5 plt.style.use('fivethirtyeight')
 6 # list of x locations for plotting
 7 x_values = list(range(len(importances)))
```

```
 8 # Make a bar chart
 9 plt.bar(x_values, importances, orientation = 'vertical')
10 # Tick labels for x axis
11 plt.xticks(x_values, feature_list, rotation='vertical')
12 # Axis labels and title
13 plt.ylabel('Importance'); plt.xlabel('Variable'); plt.title('Variable Importances');
```
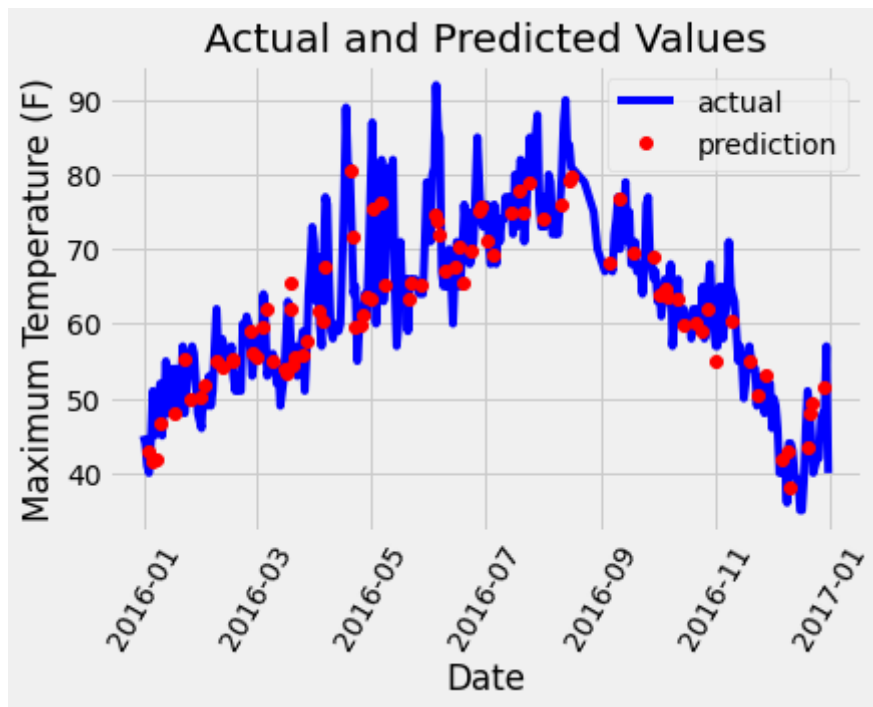


```
 1 # Use datetime for creating date objects for plotting
 2 import datetime
 3 # Dates of training values
 4 months = features[:, feature_list.index('month')]
 5 days = features[:, feature_list.index('day')]
 6 years = features[:, feature_list.index('year')]
 7 # List and then convert to datetime object
 8 dates = [str(int(year)) + '-' + str(int(month)) + '-' + str(int(day)) for year, month,
 9 dates = [datetime.datetime.strptime(date, '%Y-%m-%d') for date in dates]
10 # Dataframe with true values and dates
11 true_data = pd.DataFrame(data = {'date': dates, 'actual': labels})
12 # Dates of predictions
13 months = test_features[:, feature_list.index('month')]
14 days = test_features[:, feature_list.index('day')]
15 years = test_features[:, feature_list.index('year')]
16 # Column of dates
17 test_dates = [str(int(year)) + '-' + str(int(month)) + '-' + str(int(day)) for year, mo
18 # Convert to datetime objects
19 test_dates = [datetime.datetime.strptime(date, '%Y-%m-%d') for date in test_dates]
20 # Dataframe with predictions and dates
21 predictions_data = pd.DataFrame(data = {'date': test_dates, 'prediction': predictions})
22 # Plot the actual values
23 plt.plot(true_data['date'], true_data['actual'], 'b-', label = 'actual')
24 # Plot the predicted values
25 plt.plot(predictions_data['date'], predictions_data['prediction'], 'ro', label = 'predi
26 plt.xticks(rotation = '60');
```

```
26 plt.xticks(rotation = '60'),
27 plt.legend()
28 # Graph labels
29 plt.xlabel('Date'); plt.ylabel('Maximum Temperature (F)'); plt.title('Actual and Predic
```
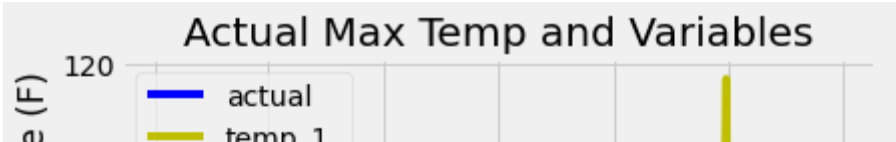


```
 1 # Make the data accessible for plotting
 2 true_data['temp_1'] = features[:, feature_list.index('temp_1')]
 3 true_data['average'] = features[:, feature_list.index('average')]
 4 true_data['friend'] = features[:, feature_list.index('friend')]
 5 # Plot all the data as lines
 6 plt.plot(true_data['date'], true_data['actual'], 'b-', label  = 'actual', alpha = 1.0)
 7 plt.plot(true_data['date'], true_data['temp_1'], 'y-', label  = 'temp_1', alpha = 1.0)
 8 plt.plot(true_data['date'], true_data['average'], 'k-', label = 'average', alpha = 0.8)
 9 plt.plot(true_data['date'], true_data['friend'], 'r-', label = 'friend', alpha = 0.3)
10 # Formatting plot
11 plt.legend(); plt.xticks(rotation = '60');
12 # Lables and title
13 plt.xlabel('Date'); plt.ylabel('Maximum Temperature (F)'); plt.title('Actual Max Temp a
14
```

1