

```
In [585]: 1 from sklearn.datasets import load_iris
2 from sklearn.tree import DecisionTreeClassifier
3
4 iris = load_iris()
5 X = iris.data[:, 2:]
6 y = iris.target_names
7
8 y
```

```
Out[585]: array(['setosa', 'versicolor', 'virginica'], dtype='<U10')
```

```
In [485]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 from scipy.stats import norm
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.decomposition import PCA
8 from scipy import stats
9 import warnings
10 import matplotlib.pyplot as plt
11
12
13 #warnings.filterwarnings('ignore')
14 %matplotlib inline
```

```
In [531]: 1 df_train = pd.read_csv('train.csv')
2 df_test = pd.read_csv('test.csv')
```

```
In [532]: 1 df_train.shape
2
```

```
Out[532]: (891, 12)
```

```
In [533]: 1 df_test.shape
```

```
Out[533]: (418, 11)
```

```
In [534]: 1 #number of rows with missing values in the dataframe!
2 df_train.isnull().values.ravel().sum()
3
```

```
Out[534]: 866
```

```
In [535]: 1 df_train.isna().sum()
```

```
Out[535]: PassengerId      0
Survived      0
Pclass        0
Name          0
Sex           0
Age          177
SibSp         0
Parch         0
Ticket        0
Fare          0
Cabin        687
Embarked      2
dtype: int64
```

```
In [536]: 1 df_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
PassengerId    891 non-null int64
Survived       891 non-null int64
Pclass         891 non-null int64
Name           891 non-null object
Sex            891 non-null object
Age           714 non-null float64
SibSp          891 non-null int64
Parch          891 non-null int64
Ticket         891 non-null object
Fare           891 non-null float64
Cabin          204 non-null object
Embarked       889 non-null object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.6+ KB
```

```
In [537]: 1 def missing_percentage(series):
2         num = series.isnull().sum()
3         den = series.count()
4         return 100*(num/den)
```

```
In [538]: 1 missing_percentage(df_train)
```

```
Out[538]: PassengerId      0.000000  
Survived      0.000000  
Pclass        0.000000  
Name          0.000000  
Sex           0.000000  
Age           24.789916  
SibSp         0.000000  
Parch         0.000000  
Ticket        0.000000  
Fare          0.000000  
Cabin         336.764706  
Embarked      0.224972  
dtype: float64
```

```
In [ ]: 1
```

```
In [539]: 1 df_train['Ticket'].value_counts()
```

```
Out[539]: 347082          7
          1601          7
          CA. 2343       7
          3101295        6
          CA 2144         6
          347088         6
          382652         5
          S.O.C. 14879    5
          W./C. 6608      4
          347077         4
          17421          4
          2666           4
          19950          4
          349909         4
          113760         4
          113781         4
          4133           4
          LINE           4
          PC 17757        4
          F.C.C. 13529    3
          C.A. 34651      3
          PC 17760        3
          248727         3
          239853         3
          35273          3
          230080         3
          13502          3
          347742         3
          29106          3
          PC 17755        3
          ..
          A/4. 34244      1
          349216          1
          PC 17595        1
          14311          1
          348121          1
          347063          1
          SOTON/O.Q. 392087 1
          SOTON/OQ 392086    1
          347468          1
          250653          1
          3474           1
          7553           1
          A/4 45380        1
          226593          1
          2648           1
          CA. 2314         1
          11753           1
          28665           1
          PC 17756         1
          113804          1
          2677            1
          2662            1
          SC/AH 29037       1
          367232           1
```

```
349212          1
C.A. 29566      1
A/5 2817        1
113050          1
350060          1
363592          1
Name: Ticket, Length: 681, dtype: int64
```

```
In [540]: 1 df_train.loc[df_train['Ticket'] == '2666']
          2
```

Out[540]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
448	449	1	3	Baclini, Miss. Marie Catherine	female	5.00	2	1	2666	19.2583	NaN
469	470	1	3	Baclini, Miss. Helene Barbara	female	0.75	2	1	2666	19.2583	NaN
644	645	1	3	Baclini, Miss. Eugenie	female	0.75	2	1	2666	19.2583	NaN
858	859	1	3	Baclini, Mrs. Solomon (Latifa Qurban)	female	24.00	0	3	2666	19.2583	NaN

```
In [ ]: 1
```

```
In [541]: 1 df_train.loc[df_train['Ticket'] == '347082']
          2
```

```
Out[541]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
13	14	0	3	Andersson, Mr. Anders Johan	male	39.0	1	5	347082	31.275	NaN
119	120	0	3	Andersson, Miss. Ellis Anna Maria	female	2.0	4	2	347082	31.275	NaN
541	542	0	3	Andersson, Miss. Ingeborg Constanzia	female	9.0	4	2	347082	31.275	NaN
542	543	0	3	Andersson, Miss. Sigrid Elisabeth	female	11.0	4	2	347082	31.275	NaN
610	611	0	3	Andersson, Mrs. Anders Johan (Alfrida Konstant...	female	39.0	1	5	347082	31.275	NaN
813	814	0	3	Andersson, Miss. Ebba Iris Alfrida	female	6.0	4	2	347082	31.275	NaN
850	851	0	3	Andersson, Master. Sigvard Harald Elias	male	4.0	4	2	347082	31.275	NaN

```
In [ ]: 1
```

```
In [542]: 1 #df_train = df_train.drop('Cabin', axis = 1)
```

```
In [543]: 1 len(df_train['Name'].str.split(',').str[0])
```

```
Out[543]: 891
```

```
In [544]: 1 len(set(df_train['Name'].str.split(',').str[0]))
```

```
Out[544]: 667
```

```
In [545]: 1 df_train['Name'] = df_train['Name'].str.split(',').str[0].value_counts()
```

```
In [546]: 1 #df_train['Name'].value_counts()
```

```
In [547]: 1 #df_train = df_train.drop(labels = 'Name', axis = 1)
```

```
In [548]: 1 #len(df_train['Cabin'])
```

```
In [549]: 1 y = df_train['Survived']  
2 y = y.astype('int64')
```

```
In [550]: 1 le = preprocessing.LabelEncoder()  
2 df_train['Sex'] = le.fit_transform(df_train.Sex.values)  
3 df_train['Ticket'] = le.fit_transform(df_train.Ticket.values)  
4 df_train['Name'] = le.fit_transform(df_train.Name.values)  
5 df_train['Embarked'] = le.fit_transform(df_train['Embarked']).astype(str)
```

```
In [551]: 1 df_train = df_train.drop(labels = 'Cabin', axis = 1 )
```

```
In [552]: 1 x = df_train.drop(labels = 'Survived', axis = 1 )
```

```
In [555]: 1 from sklearn.impute import SimpleImputer  
2 imp = SimpleImputer(missing_values=np.nan, strategy='mean')  
3 x = imp.fit_transform(x)  
4 #x
```

```
In [556]: 1 from sklearn.decomposition import PCA  
2 import numpy as np  
3  
4 pca = PCA(n_components=2)  
5 x_transformed = pca.fit_transform(x)  
6
```

```
In [557]: 1 print(pca.explained_variance_ratio_)  
  
[0.45487329 0.27775175]
```

```
In [558]: 1 y = df_train['Survived']  
2 set(y.values)
```

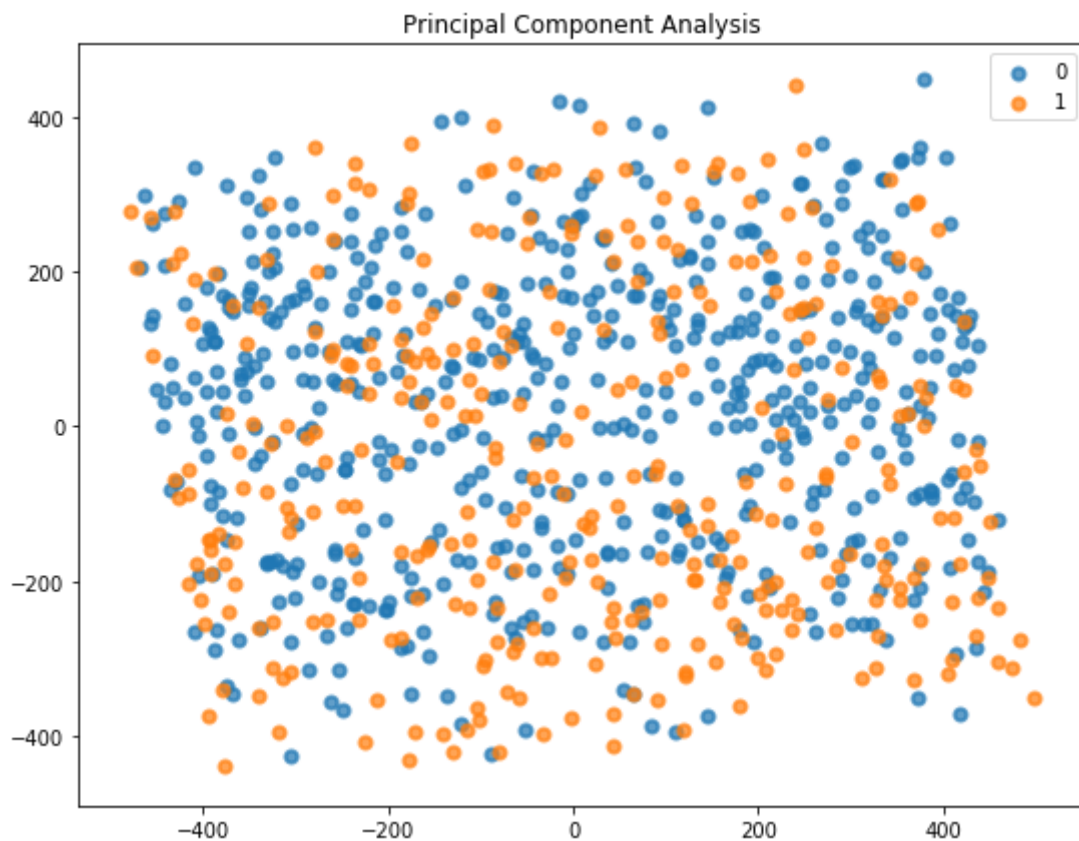
```
Out[558]: {0, 1}
```

```
In [559]: 1 x_transformed.shape
```

```
Out[559]: (891, 2)
```

```
In [560]: 1 import numpy as np  
2 import matplotlib.pyplot as plt  
3 from mpl_toolkits.mplot3d import Axes3D  
4
```

```
In [561]: 1
2 plt.figure(figsize=(9,7))
3 lw = 2
4
5 for i, target_name in zip([0, 1], np.unique(y)):
6     plt.scatter(x_transformed[y == i, 0], x_transformed[y == i, 1], alpha=0.5,
7                 label=target_name)
8 plt.legend(loc = 'best', shadow = False, scatterpoints = 1)
9 plt.title('Principal Component Analysis')
10 plt.show()
11
```



```
In [562]: 1 #from sklearn.manifold import TSNE
```

```
In [564]: 1 #result=pd.DataFrame(x_transformed, columns=['PCA%i' % i for i in range
2 #result
```

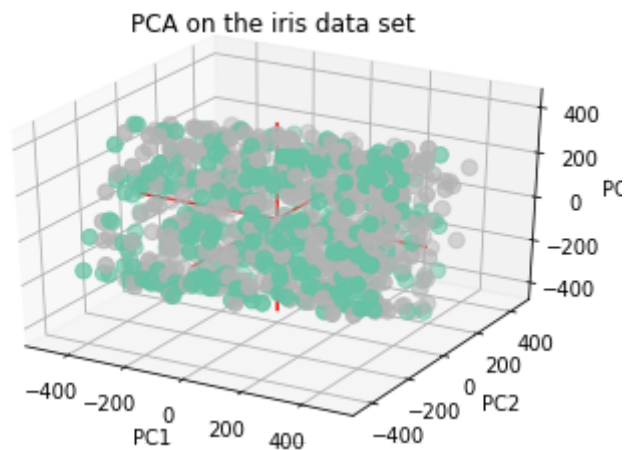


```

In [565]: 1 # Plot initialisation
2 fig = plt.figure()
3 ax = fig.add_subplot(111, projection='3d')
4 ax.scatter(result['PCA0'], result['PCA1'], result['PCA2'], c=my_color,
5
6 # make simple, bare axis lines through space:
7 xAxisLine = ((min(result['PCA0']), max(result['PCA0'])), (0, 0), (0,0))
8 ax.plot(xAxisLine[0], xAxisLine[1], xAxisLine[2], 'r')
9 yAxisLine = ((0, 0), (min(result['PCA1']), max(result['PCA1'])), (0,0))
10 ax.plot(yAxisLine[0], yAxisLine[1], yAxisLine[2], 'r')
11 zAxisLine = ((0, 0), (0,0), (min(result['PCA2']), max(result['PCA2'])))
12 ax.plot(zAxisLine[0], zAxisLine[1], zAxisLine[2], 'r')
13
14 # label the axes
15 ax.set_xlabel("PC1")
16 ax.set_ylabel("PC2")
17 ax.set_zlabel("PC3")
18 ax.set_title("PCA on the iris data set")
19 #plt.show()
20

```

Out[565]: Text(0.5, 0.92, 'PCA on the iris data set')



```

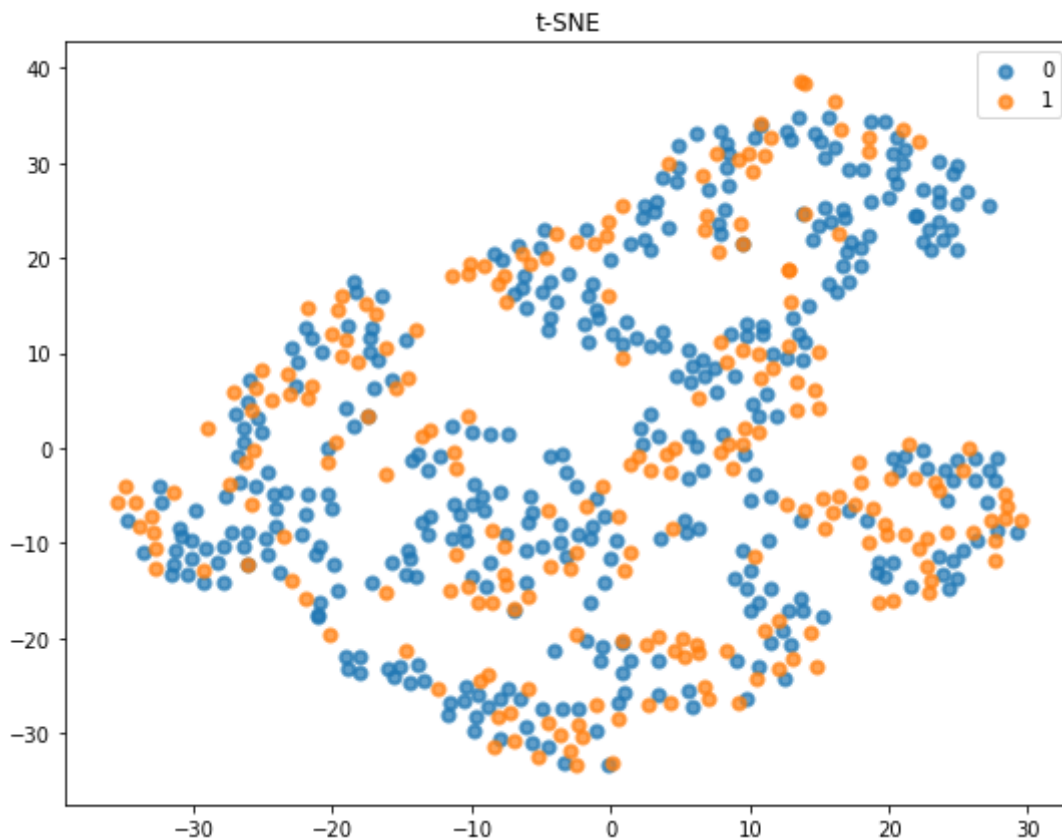
In [566]: 1 from sklearn import tree
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.3
4

```

```
In [567]: 1 from sklearn.manifold import TSNE
2 tsne = TSNE(n_components=2)
3 X_embedded = tsne.fit_transform(X_train)
4 X_embedded
```

```
Out[567]: array([[ -6.9513154, -17.071068 ],
 [  9.500626 ,  21.470192 ],
 [  8.347942 ,  32.102566 ],
 ...,
 [ 20.496828 ,  27.774937 ],
 [ 10.424785 , -24.273827 ],
 [-13.52143  ,  -7.786212 ]], dtype=float32)
```

```
In [568]: 1 plt.figure(figsize=(9,7))
2 lw = 2
3
4 for i, target_name in zip([0, 1], np.unique(y_train)):
5     plt.scatter(X_embedded[y_train == i, 0], X_embedded[y_train == i, 1],
6               label=target_name)
7 plt.legend(loc = 'best', shadow = False, scatterpoints = 1)
8 plt.title('t-SNE')
9 plt.show()
10
```



```
In [569]: 1 clf = tree.DecisionTreeClassifier(max_depth =3)
2 clf = clf.fit(X_train, y_train)
```

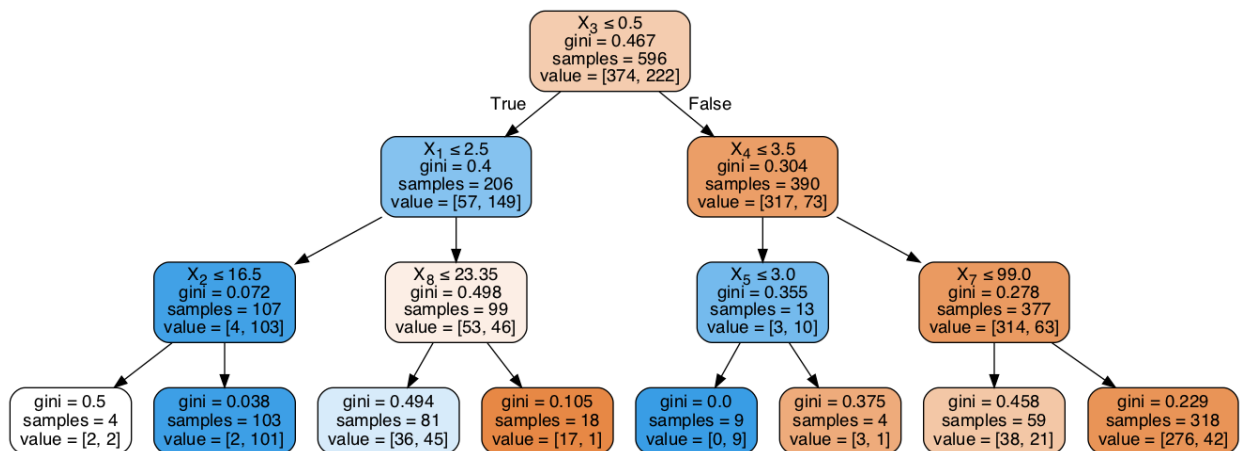
```
In [570]: 1 y_pred = clf.predict(X_test)
```

```
In [571]: 1 from sklearn.metrics import accuracy_score
2 accuracy_score(y_test, y_pred)
3
```

Out[571]: 0.8203389830508474

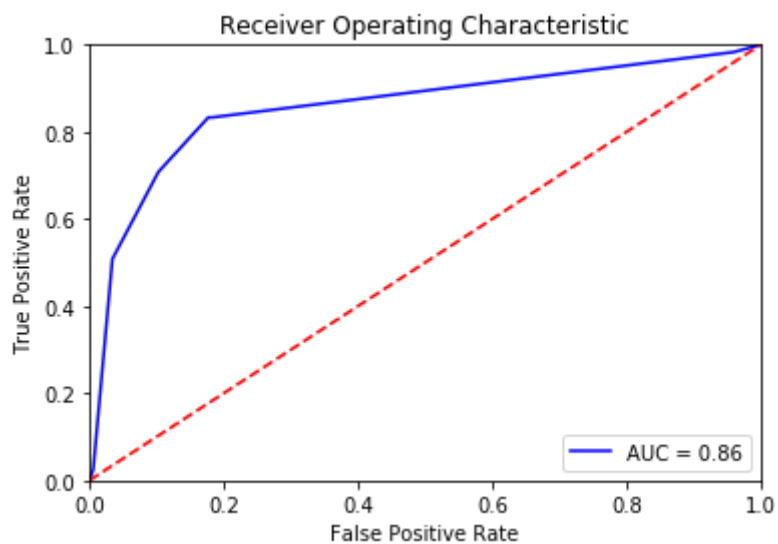
```
In [572]: 1 from sklearn.externals.six import StringIO
2 from IPython.display import Image
3 from sklearn.tree import export_graphviz
4 import pydotplus
5 dot_data = StringIO()
6 export_graphviz(clf, out_file=dot_data,
7                 filled=True, rounded=True,
8                 special_characters=True)
9 graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
10 Image(graph.create_png())
```

Out[572]:

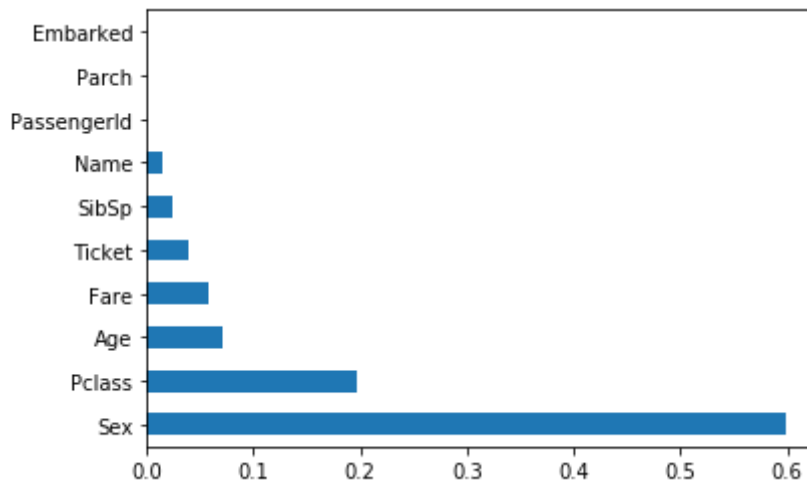


```
In [438]: 1 import sklearn.metrics as metrics
2 # calculate the fpr and tpr for all thresholds of the classification
3 probs = clf.predict_proba(X_test)
4 preds = probs[:,1]
5 fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
6 roc_auc = metrics.auc(fpr, tpr)
7
8
```

```
In [439]: 1 # method 1: plt
2 import matplotlib.pyplot as plt
3 plt.title('Receiver Operating Characteristic')
4 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
5 plt.legend(loc = 'lower right')
6 plt.plot([0, 1], [0, 1], 'r--')
7 plt.xlim([0, 1])
8 plt.ylim([0, 1])
9 plt.ylabel('True Positive Rate')
10 plt.xlabel('False Positive Rate')
11 plt.show()
```



```
In [404]: 1 feat_importances = pd.Series(clf.feature_importances_, index=df_train.dr
2 feat_importances.nlargest(25).plot(kind='barh')
3 plt.show()
```



```
In [415]: 1 from sklearn.ensemble import ExtraTreesClassifier, GradientBoostingClassifier
2 model = ExtraTreesClassifier()
3 model.fit(X_train, y_train)
4
5
6
```

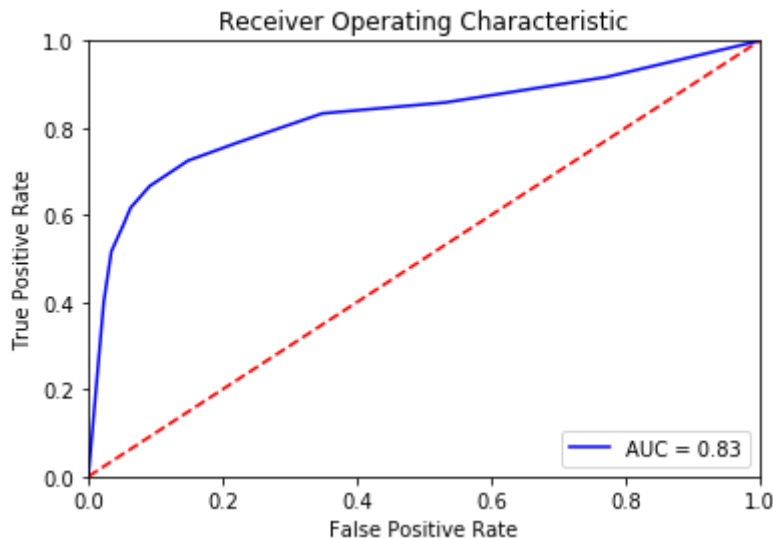
```
Out[415]: ExtraTreesClassifier(bootstrap=False, class_weight=None, criterion='gini',
                                max_depth=None, max_features='auto', max_leaf_nodes=
                                None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=1, min_samples_split=2,
                                min_weight_fraction_leaf=0.0, n_estimators=10, n_job
                                s=None,
                                oob_score=False, random_state=None, verbose=0,
                                warm_start=False)
```

```
In [450]: 1 #probs#[ :,1]
```

```

In [446]: 1 import sklearn.metrics as metrics
2 # calculate the fpr and tpr for all thresholds of the classification
3 probs = model.predict_proba(X_test)
4 #preds = probs[:,1]
5 fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
6 roc_auc = metrics.auc(fpr, tpr)
7
8
9 # method 1: plt
10 plt.title('Receiver Operating Characteristic')
11 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
12 plt.legend(loc = 'lower right')
13 plt.plot([0, 1], [0, 1], 'r--')
14 plt.xlim([0, 1])
15 plt.ylim([0, 1])
16 plt.ylabel('True Positive Rate')
17 plt.xlabel('False Positive Rate')
18 plt.show()

```

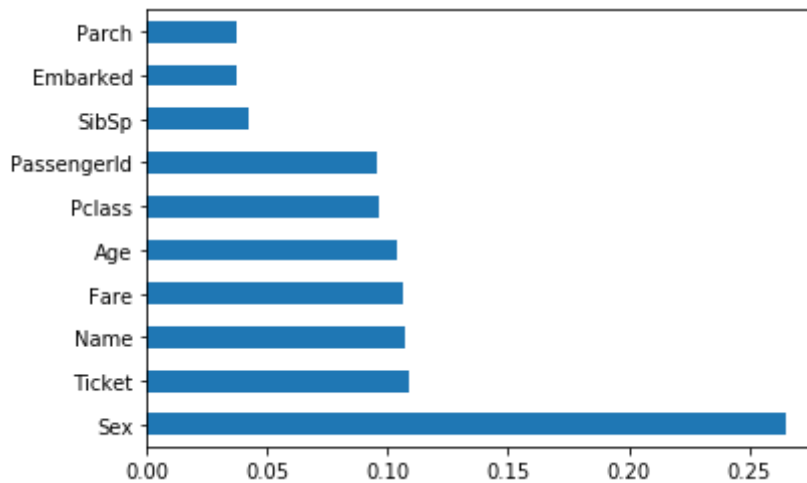


```

In [416]: 1 feat_importances = pd.Series(model.feature_importances_, index=df_train.
2 #len(df_train.columns)
3 #model.feature_importances_
4 #df_train.drop(labels= 'Survived', axis = 1).columns
5

```

```
In [417]: 1 #feat_importances = pd.Series(model.feature_importances_, index=df_train
2 feat_importances.nlargest(25).plot(kind='barh')
3 plt.show()
```



```
In [463]: 1 y_pred = model.predict(X_test)
2 accuracy_score(y_test, y_pred)
3
```

Out[463]: 0.8101694915254237

```
In [454]: 1 modelgc = GradientBoostingClassifier()
2 modelgc.fit(X_train, y_train)
3
```

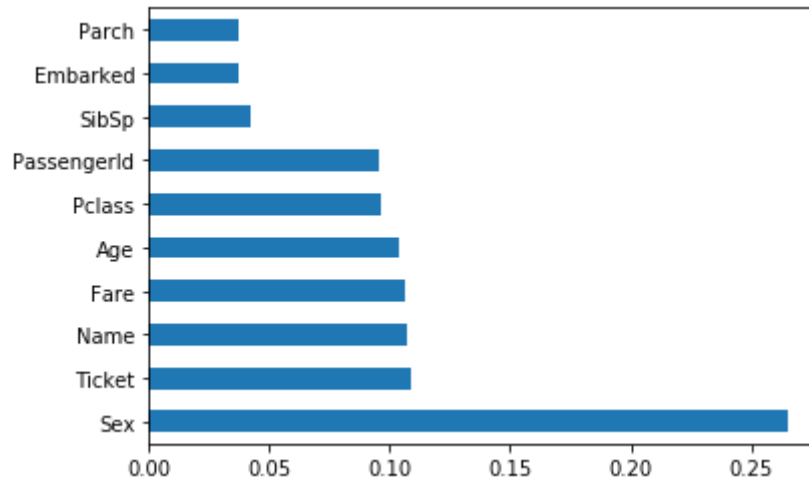
Out[454]: GradientBoostingClassifier(criterion='friedman_mse', init=None, learning_rate=0.1, loss='deviance', max_depth=3, max_features=None, max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=10, n_iter_no_change=None, presort='auto', random_state=None, subsample=1.0, tol=0.0001, validation_fraction=0.1, verbose=0, warm_start=False)

```
In [482]: 1 y_predgc = modelgc.predict(X_test)
          2 accuracy_score(y_test, y_pred)
          3
```

Out[482]: 0.8101694915254237

```
In [484]: 1 #feat_importances = pd.Series(modelgc.feature_importances_,index=df_tra
          2 #df_train
```

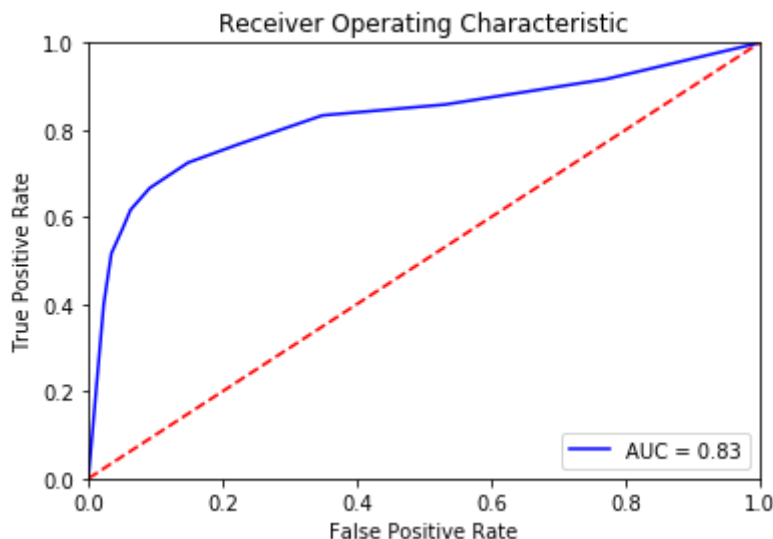
```
In [469]: 1 feat_importances.nlargest(25).plot(kind='barh')
          2 plt.show()
```




```

In [470]: 1 import sklearn.metrics as metrics
2 # calculate the fpr and tpr for all thresholds of the classification
3 probs = modelgc.predict_proba(X_test)
4 #preds = probs[:,1]
5 fpr, tpr, threshold = metrics.roc_curve(y_test, probs)
6 roc_auc = metrics.auc(fpr, tpr)
7
8
9 # method 1: plt
10 plt.title('Receiver Operating Characteristic')
11 plt.plot(fpr, tpr, 'b', label = 'AUC = %0.2f' % roc_auc)
12 plt.legend(loc = 'lower right')
13 plt.plot([0, 1], [0, 1], 'r--')
14 plt.xlim([0, 1])
15 plt.ylim([0, 1])
16 plt.ylabel('True Positive Rate')
17 plt.xlabel('False Positive Rate')
18 plt.show()

```



```

In [472]: 1 y_predgc

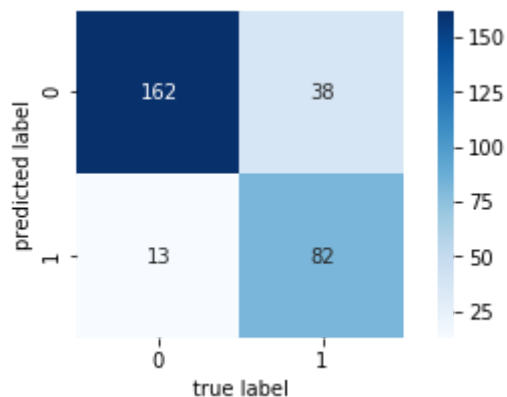
```

```

Out[472]: array([0, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 0, 0,
1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0,
1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1,
0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,
1, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0,
0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0,
0, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0,
0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0,
1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0,
0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 1, 0, 1, 0,
0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0,
1, 0, 1, 1, 0, 0, 1, 1, 0])

```

```
In [481]: 1 from sklearn.metrics import confusion_matrix, accuracy_score
2
3 mat = confusion_matrix(y_test, y_predgc)
4
5 plt.figure(figsize=(6, 3))
6
7 sns.heatmap(mat.T, xticklabels=np.unique(y_train),
8             annot=True, fmt="d", square=True, yticklabels=np.unique(y_train),
9             plt.xlabel('true label')
10            plt.ylabel('predicted label');
11
```



```
In [580]: 1 from sklearn.metrics import precision_recall_curve
2 from sklearn.metrics import plot_precision_recall_curve
3 import matplotlib.pyplot as plt
4
5 disp = plot_precision_recall_curve(modelgc, X_test, y_test)
6 disp.ax_.set_title('2-class Precision-Recall curve: AP={0:0.2f}'.format(
```

```
-----
--
ImportError                                Traceback (most recent call last)
```

```
<ipython-input-580-b73fd41f808d> in <module>
      1 from sklearn.metrics import precision_recall_curve
----> 2 from sklearn.metrics import plot_precision_recall_curve
      3 import matplotlib.pyplot as plt
      4
      5 disp = plot_precision_recall_curve(modelgc, X_test, y_test)
```

```
ImportError: cannot import name 'plot_precision_recall_curve' from 'sklearn.metrics' (/Users/shradhitsu/anaconda3/envs/python37charm/lib/python3.7/site-packages/sklearn/metrics/__init__.py)
```

```
In [578]: 1 from sklearn import svm, datasets
2 from sklearn.model_selection import train_test_split
3 import numpy as np
4 random_state = np.random.RandomState(0)
5
6 # Create a simple classifier
7 classifier = svm.LinearSVC(random_state=random_state)
8 classifier.fit(X_train, y_train)
9 y_score = classifier.decision_function(X_test)
```

```
In [581]: 1 from sklearn.metrics import average_precision_score
2 average_precision = average_precision_score(y_test, y_score)
3
4 print('Average precision-recall score: {0:0.2f}'.format(
5     average_precision))
6
```

Average precision-recall score: 0.65

```
In [582]: 1 from sklearn.metrics import precision_recall_curve
2 from sklearn.metrics import plot_precision_recall_curve
3 import matplotlib.pyplot as plt
4
5 disp = plot_precision_recall_curve(modelgc, X_test, y_test)
6 disp.ax_.set_title('2-class Precision-Recall curve: AP={0:0.2f}'.format
```

```
-----
--
ImportError                                Traceback (most recent call las
t)
<ipython-input-582-b73fd41f808d> in <module>
      1 from sklearn.metrics import precision_recall_curve
----> 2 from sklearn.metrics import plot_precision_recall_curve
      3 import matplotlib.pyplot as plt
      4
      5 disp = plot_precision_recall_curve(modelgc, X_test, y_test)

ImportError: cannot import name 'plot_precision_recall_curve' from 'sklea
rn.metrics' (/Users/shradhitsuBudhi/anaconda3/envs/python37charm/lib/pyth
on3.7/site-packages/sklearn/metrics/__init__.py)
```

```
In [ ]:
```

```
1
```