```python
In [1]:  import pandas as pd
         import numpy as np
         import matplotlib.pyplot as plt
         %matplotlib inline
         diabetes = pd.read_csv('diabetes.csv')
         print(diabetes.columns)
```
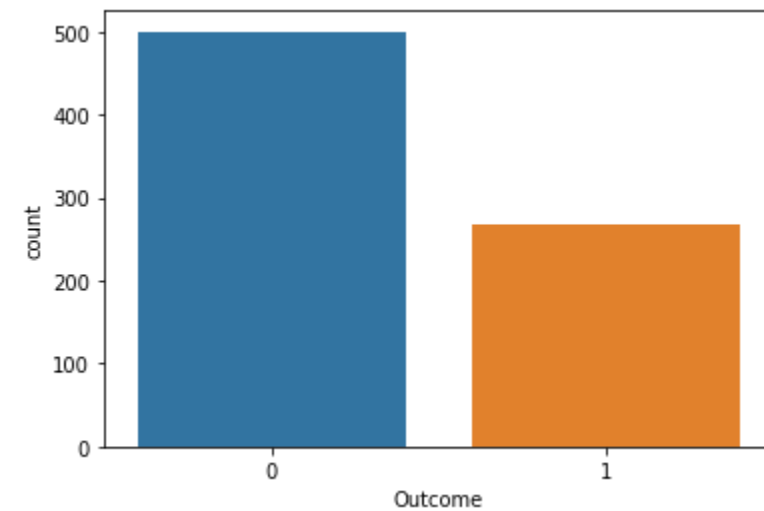
```
Index(['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin',
       'BMI', 'DiabetesPedigreeFunction', 'Age', 'Outcome'],
      dtype='object')
```

```python
In [2]:  print("dimension of diabetes data: {}".format(diabetes.shape))
```
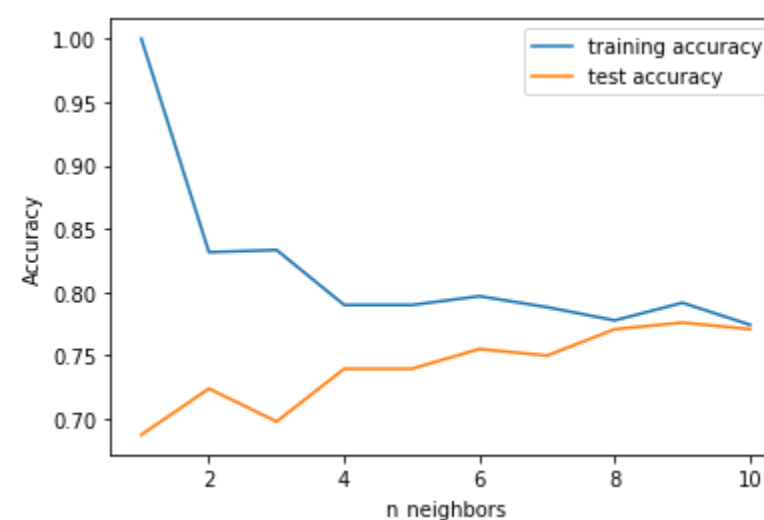
```
dimension of diabetes data: (768, 9)
```

```python
In [3]:  import seaborn as sns
         sns.countplot(diabetes['Outcome'],label="Count")
```

```
Out[3]:  <matplotlib.axes._subplots.AxesSubplot at 0x175b60fa948>
```



```python
In [4]:  from sklearn.model_selection import train_test_split
         X_train, X_test, y_train, y_test = train_test_split(diabetes.loc[:, diabetes.columns != 'Out
         come'], diabetes['Outcome'], stratify=diabetes['Outcome'], random_state=66)
         from sklearn.neighbors import KNeighborsClassifier
         training_accuracy = []
         test_accuracy = []
         # try n_neighbors from 1 to 10
         neighbors_settings = range(1, 11)
         for n_neighbors in neighbors_settings:
             # build the model
             knn = KNeighborsClassifier(n_neighbors=n_neighbors)
             knn.fit(X_train, y_train)
             # record training set accuracy
             training_accuracy.append(knn.score(X_train, y_train))
             # record test set accuracy
             test_accuracy.append(knn.score(X_test, y_test))
         plt.plot(neighbors_settings, training_accuracy, label="training accuracy")
         plt.plot(neighbors_settings, test_accuracy, label="test accuracy")
         plt.ylabel("Accuracy")
         plt.xlabel("n_neighbors")
         plt.legend()
         plt.savefig('knn_compare_model')
```



```python
In [5]:  knn = KNeighborsClassifier(n_neighbors=9)
         knn.fit(X_train, y_train)
         print('Accuracy of K-NN classifier on training set: {:.2f}'.format(knn.score(X_train, y_trai
         n)))
         print('Accuracy of K-NN classifier on test set: {:.2f}'.format(knn.score(X_test, y_test)))
```

```
Accuracy of K-NN classifier on training set: 0.79
Accuracy of K-NN classifier on test set: 0.78
```

```python
In [6]:  from sklearn.linear_model import LogisticRegression
         logreg = LogisticRegression().fit(X_train, y_train)
         print("Training set score: {:.3f}".format(logreg.score(X_train, y_train)))
         print("Test set score: {:.3f}".format(logreg.score(X_test, y_test)))
```

```
Training set score: 0.781
Test set score: 0.771
```

```
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:432: FutureWarnin
g: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warnin
g.
  FutureWarning)
```

```python
In [9]:  tree = DecisionTreeClassifier(max_depth=3, random_state=0)
         tree.fit(X_train, y_train)
         print("Accuracy on training set: {:.3f}".format(tree.score(X_train, y_train)))
         print("Accuracy on test set: {:.3f}".format(tree.score(X_test, y_test)))
```

```
Accuracy on training set: 0.773
Accuracy on test set: 0.740
```

```python
In [10]: from sklearn.ensemble import RandomForestClassifier
         rf = RandomForestClassifier(n_estimators=100, random_state=0)
         rf.fit(X_train, y_train)
         print("Accuracy on training set: {:.3f}".format(rf.score(X_train, y_train)))
         print("Accuracy on test set: {:.3f}".format(rf.score(X_test, y_test)))
```

```
Accuracy on training set: 1.000
Accuracy on test set: 0.786
```

```python
In [11]: rf1 = RandomForestClassifier(max_depth=3, n_estimators=100, random_state=0)
         rf1.fit(X_train, y_train)
         print("Accuracy on training set: {:.3f}".format(rf1.score(X_train, y_train)))
         print("Accuracy on test set: {:.3f}".format(rf1.score(X_test, y_test)))
```

```
Accuracy on training set: 0.800
Accuracy on test set: 0.755
```