Name: Apoorva Dhalpawar
Roll No.: 17CO012
Subject: Machine Learning
Class: BE Computer 1

Assignment 2:
Problem Statement:

Search a Medical related dataset
1. Download the dataset
2. Perform pre-processing on the dataset
3. Use this dataset to build a Naive Bayes Classifier
4. Use this dataset to build a Decision Tree Classifier
5. Compare the results and comment

Tool used:
Jupyter Notebook(Python).

Dataset:
Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image. The class has values that classify the tumor into "Malign" or "Benign"

Implementation:
Below attached is the converted notebook file with output for both the classifiers.

Conclusion:
Test data set performance of both the models are:
　　1. Decision Tree: 95.8%
　　2. Naive Bayes: 95.3%

In [14]:

```python
from sklearn.datasets import load_breast_cancer
from sklearn.tree import DecisionTreeClassifier        #Decision Tree
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.tree import export_graphviz
import matplotlib.pylab as plt
from sklearn.metrics import classification_report
from sklearn import metrics
from sklearn.naive_bayes import GaussianNB             #Naive Bayes
import numpy as np
import graphviz
%matplotlib inline
```

In [15]:

```python
#load the breast cancer data and few EDA
cancer = load_breast_cancer()
print(cancer.DESCR)
```

```
.. _breast_cancer_dataset:
```

Breast cancer wisconsin (diagnostic) dataset
--------------------------------------------

**Data Set Characteristics:**

    :Number of Instances: 569

    :Number of Attributes: 30 numeric, predictive attributes and the class

    :Attribute Information:
        - radius (mean of distances from center to points on the perimete
r)
        - texture (standard deviation of gray-scale values)
        - perimeter
        - area
        - smoothness (local variation in radius lengths)
        - compactness (perimeter^2 / area - 1.0)
        - concavity (severity of concave portions of the contour)
        - concave points (number of concave portions of the contour)
        - symmetry
        - fractal dimension ("coastline approximation" - 1)

        The mean, standard error, and "worst" or largest (mean of the thre
e
        largest values) of these features were computed for each image,
        resulting in 30 features.  For instance, field 3 is Mean Radius, f
ield
        13 is Radius SE, field 23 is Worst Radius.

        - class:
                - WDBC-Malignant
                - WDBC-Benign

    :Summary Statistics:

    ===================================== ====== ======
                                           Min    Max
    ===================================== ====== ======
    radius (mean):                        6.981  28.11
    texture (mean):                       9.71   39.28
    perimeter (mean):                     43.79  188.5
    area (mean):                          143.5  2501.0
    smoothness (mean):                    0.053  0.163
    compactness (mean):                   0.019  0.345
    concavity (mean):                     0.0    0.427
    concave points (mean):                0.0    0.201
    symmetry (mean):                      0.106  0.304
    fractal dimension (mean):             0.05   0.097
    radius (standard error):              0.112  2.873
    texture (standard error):             0.36   4.885
    perimeter (standard error):           0.757  21.98
    area (standard error):                6.802  542.2
    smoothness (standard error):          0.002  0.031
    compactness (standard error):         0.002  0.135
    concavity (standard error):           0.0    0.396
    concave points (standard error):      0.0    0.053
    symmetry (standard error):            0.008  0.079
    fractal dimension (standard error):   0.001  0.03
    radius (worst):                       7.93   36.04
```

```
    texture (worst):                      12.02  49.54
    perimeter (worst):                    50.41  251.2
    area (worst):                         185.2  4254.0
    smoothness (worst):                   0.071  0.223
    compactness (worst):                  0.027  1.058
    concavity (worst):                    0.0    1.252
    concave points (worst):               0.0    0.291
    symmetry (worst):                     0.156  0.664
    fractal dimension (worst):            0.055  0.208
    ==================================== ====== ======
```

:Missing Attribute Values: None

:Class Distribution: 212 - Malignant, 357 - Benign

:Creator:  Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian

:Donor: Nick Street

:Date: November, 1995

This is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets.
https://goo.gl/U2Uwz2

Features are computed from a digitized image of a fine needle
aspirate (FNA) of a breast mass.  They describe
characteristics of the cell nuclei present in the image.

Separating plane described above was obtained using
Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree
Construction Via Linear Programming." Proceedings of the 4th
Midwest Artificial Intelligence and Cognitive Science Society,
pp. 97-101, 1992], a classification method which uses linear
programming to construct a decision tree.  Relevant features
were selected using an exhaustive search in the space of 1-4
features and 1-3 separating planes.

The actual linear program used to obtain the separating plane
in the 3-dimensional space is that described in:
[K. P. Bennett and O. L. Mangasarian: "Robust Linear
Programming Discrimination of Two Linearly Inseparable Sets",
Optimization Methods and Software 1, 1992, 23-34].

This database is also available through the UW CS ftp server:

ftp ftp.cs.wisc.edu
cd math-prog/cpo-dataset/machine-learn/WDBC/

.. topic:: References

    - W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extra
ction
      for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on
      Electronic Imaging: Science and Technology, volume 1905, pages 861-87
0,
      San Jose, CA, 1993.
    - O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnos
is and
      prognosis via linear programming. Operations Research, 43(4), pages 5
70-577,
      July-August 1995.

- W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning tec
hniques
    to diagnose breast cancer from fine-needle aspirates. Cancer Letters
77 (1994)
    163-171.

In [16]:

```python
print(cancer.feature_names)
```

```
['mean radius' 'mean texture' 'mean perimeter' 'mean area'
 'mean smoothness' 'mean compactness' 'mean concavity'
 'mean concave points' 'mean symmetry' 'mean fractal dimension'
 'radius error' 'texture error' 'perimeter error' 'area error'
 'smoothness error' 'compactness error' 'concavity error'
 'concave points error' 'symmetry error' 'fractal dimension error'
 'worst radius' 'worst texture' 'worst perimeter' 'worst area'
 'worst smoothness' 'worst compactness' 'worst concavity'
 'worst concave points' 'worst symmetry' 'worst fractal dimension']
```

In [17]:

```python
print(cancer.target_names)
```

```
['malignant' 'benign']
```

In [18]:

```python
cancer.data
```

Out[18]:

```
array([[1.799e+01, 1.038e+01, 1.228e+02, ..., 2.654e-01, 4.601e-01,
        1.189e-01],
       [2.057e+01, 1.777e+01, 1.329e+02, ..., 1.860e-01, 2.750e-01,
        8.902e-02],
       [1.969e+01, 2.125e+01, 1.300e+02, ..., 2.430e-01, 3.613e-01,
        8.758e-02],
       ...,
       [1.660e+01, 2.808e+01, 1.083e+02, ..., 1.418e-01, 2.218e-01,
        7.820e-02],
       [2.060e+01, 2.933e+01, 1.401e+02, ..., 2.650e-01, 4.087e-01,
        1.240e-01],
       [7.760e+00, 2.454e+01, 4.792e+01, ..., 0.000e+00, 2.871e-01,
        7.039e-02]])
```

In [19]:

```python
type(cancer.data)
cancer.data.shape
```

Out[19]:

```
(569, 30)
```

## DECISION TREE

In [20]:

```python
X_train, X_test, y_train, y_test = train_test_split(cancer.data, cancer.target, random_
state=42)

training_accuracy = []
test_accuracy = []

max_dep = range(1,15)

for md in max_dep:
    tree = DecisionTreeClassifier(max_depth=md,random_state=0)
    tree.fit(X_train,y_train)
    training_accuracy.append(tree.score(X_train, y_train))
    test_accuracy.append(tree.score(X_test, y_test))

plt.plot(max_dep,training_accuracy, label='Accuracy of the training set')
plt.ylabel('Accuracy')
plt.xlabel('Max Depth')
plt.legend()

# By having larger max_depth (>5), we overfit the model into training data, so the accu
racy for training set become
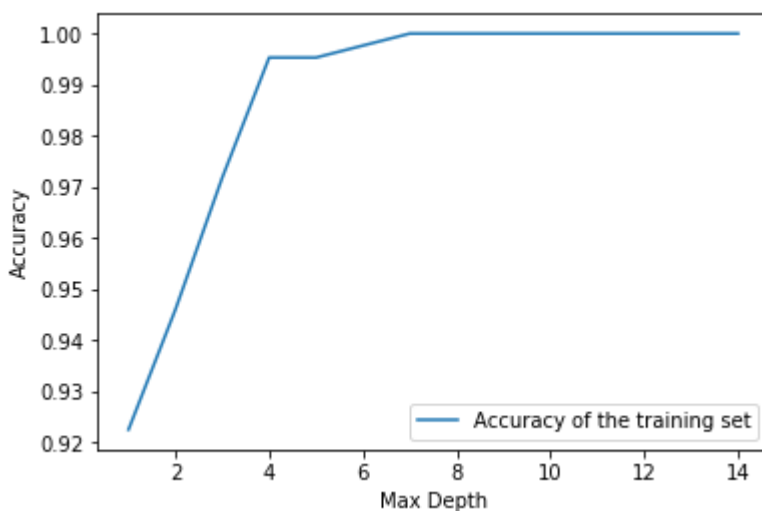# but the accuracy for test set decrease

# other parameters than can work with:
# - min_samples_leaf, max_sample_leaf
# - max_leaf_node

# by looking at plot, best result accurs when max_depth is 3
```

Out[20]:

<matplotlib.legend.Legend at 0x7f4c7aec48d0>

In [21]:

```
tree = DecisionTreeClassifier(max_depth=3,random_state=0)
tree.fit(X_train,y_train)
```

Out[21]:

```
DecisionTreeClassifier(class_weight=None, criterion='gini', 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, presort=False, random_state=0,
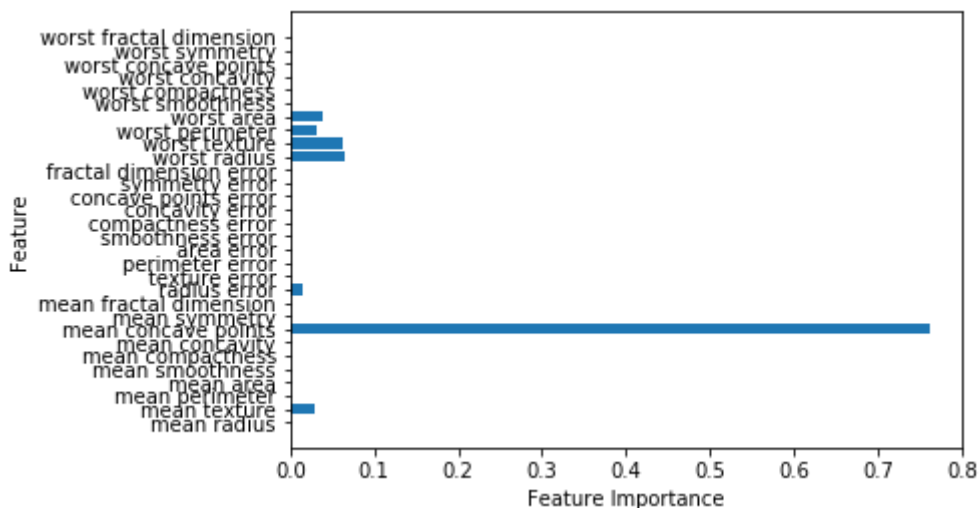            splitter='best')
```

In [22]:

```
print('Training accuracy : {}').format(tree.score(X_train, y_train))
print('Testing accuracy : {}').format(tree.score(X_test, y_test))
```

```
Training accuracy : 0.971830985915
Testing accuracy : 0.958041958042
```

In [23]:

```
#Feature Importance
n_feature = cancer.data.shape[1]
plt.barh(range(n_feature), tree.feature_importances_, align='center')
plt.yticks(np.arange(n_feature), cancer.feature_names)
plt.xlabel('Feature Importance')
plt.ylabel('Feature')
plt.show()
```



## Naive Bayes

In [24]:

```python
X_train, X_test, y_train, y_test = train_test_split(cancer.data, cancer.target, test_si
ze=0.3, random_state=109)

classifier = GaussianNB()
naive_bayes_model = classifier.fit(X_train, y_train)
y_true, y_pred = y_test, naive_bayes_model.predict(X_test)


print(classification_report(y_true, y_pred))
print("Accuracy:", metrics.accuracy_score(y_test, y_pred))
print("Precision:",metrics.precision_score(y_test, y_pred))
print("Recall:",metrics.recall_score(y_test, y_pred))

print(metrics.confusion_matrix(y_test, y_pred))
```

```
              precision    recall  f1-score   support

           0       0.97      0.90      0.93        63
           1       0.95      0.98      0.96       108

   micro avg       0.95      0.95      0.95       171
   macro avg       0.96      0.94      0.95       171
weighted avg       0.95      0.95      0.95       171

('Accuracy:', 0.9532163742690059)
('Precision:', 0.9464285714285714)
('Recall:', 0.9814814814814815)
[[ 57    6]
 [  2 106]]
```

In [ ]: