# Principal Component Analysis

In [1]:
```python
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
import seaborn as sns
%matplotlib inline
```

In [2]:
```python
from sklearn.datasets import load_breast_cancer
```

In [3]:
```python
cancer = load_breast_cancer()
```

In [4]:
```python
type(cancer)
```

Out[4]:
```
sklearn.utils.Bunch
```

In [5]:
```python
cancer.keys()
```

Out[5]:
```
dict_keys(['data', 'target', 'frame', 'target_names', 'DESCR', 'feature_names', 'filename', 'data_module'])
```

In [6]:
```python
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 perimeter)
        - 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 three
        worst/largest values) of these features were computed for each image,
        resulting in 30 features.  For instance, field 0 is Mean Radius, field
        10 is Radius SE, field 20 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 extraction
      for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on
      Electronic Imaging: Science and Technology, volume 1905, pages 861-870,
      San Jose, CA, 1993.
    - O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and
      prognosis via linear programming. Operations Research, 43(4), pages 570-577,
      July-August 1995.
    - W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques
      to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994)
      163-171.

```
In [8]:   df = pd.DataFrame(cancer['data'],columns=cancer['feature_names'])
```

```
In [10]:  #So it is hard to visualize a high dimensional data especially here where we have
          # 30 numeric variables (attributes) hence our data has 30 dimensions.

          #So let's use PCA to find the main 2 components and transform it into 2D
```

```
In [12]:  from sklearn.preprocessing import StandardScaler
```

```
In [13]:  scaler = StandardScaler()
```

```
In [14]:  scaler.fit(df)
```

```
Out[14]:  StandardScaler()
```

```
In [15]:  scaled_data = scaler.transform(df)
```

```
In [16]:  #PCA
          from sklearn.decomposition import PCA
```

```
In [17]:  pca = PCA(n_components=2)
```

```
In [18]:  pca.fit(scaled_data)
```

```
Out[18]:  PCA(n_components=2)
```

```
In [19]:  x_pca = pca.transform(scaled_data)
```

```
In [20]:  scaled_data.shape
```
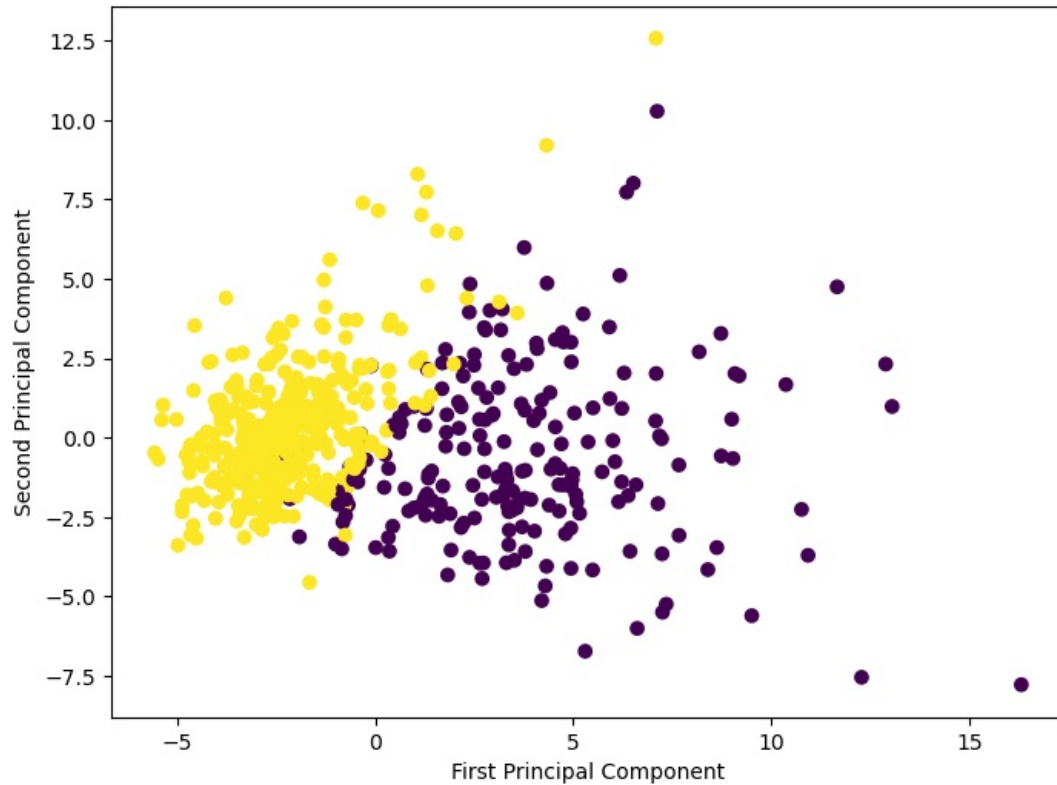
```
Out[20]:  (569, 30)
```

```
In [21]:  x_pca.shape
```

```
Out[21]:  (569, 2)
```

```
In [24]:  plt.figure(figsize=(8,6))
          plt.scatter(x_pca[:,0],x_pca[:,1],c=cancer['target'])
          plt.xlabel('First Principal Component')
          plt.ylabel('Second Principal Component')
```

```
Out[24]:  Text(0, 0.5, 'Second Principal Component')
```
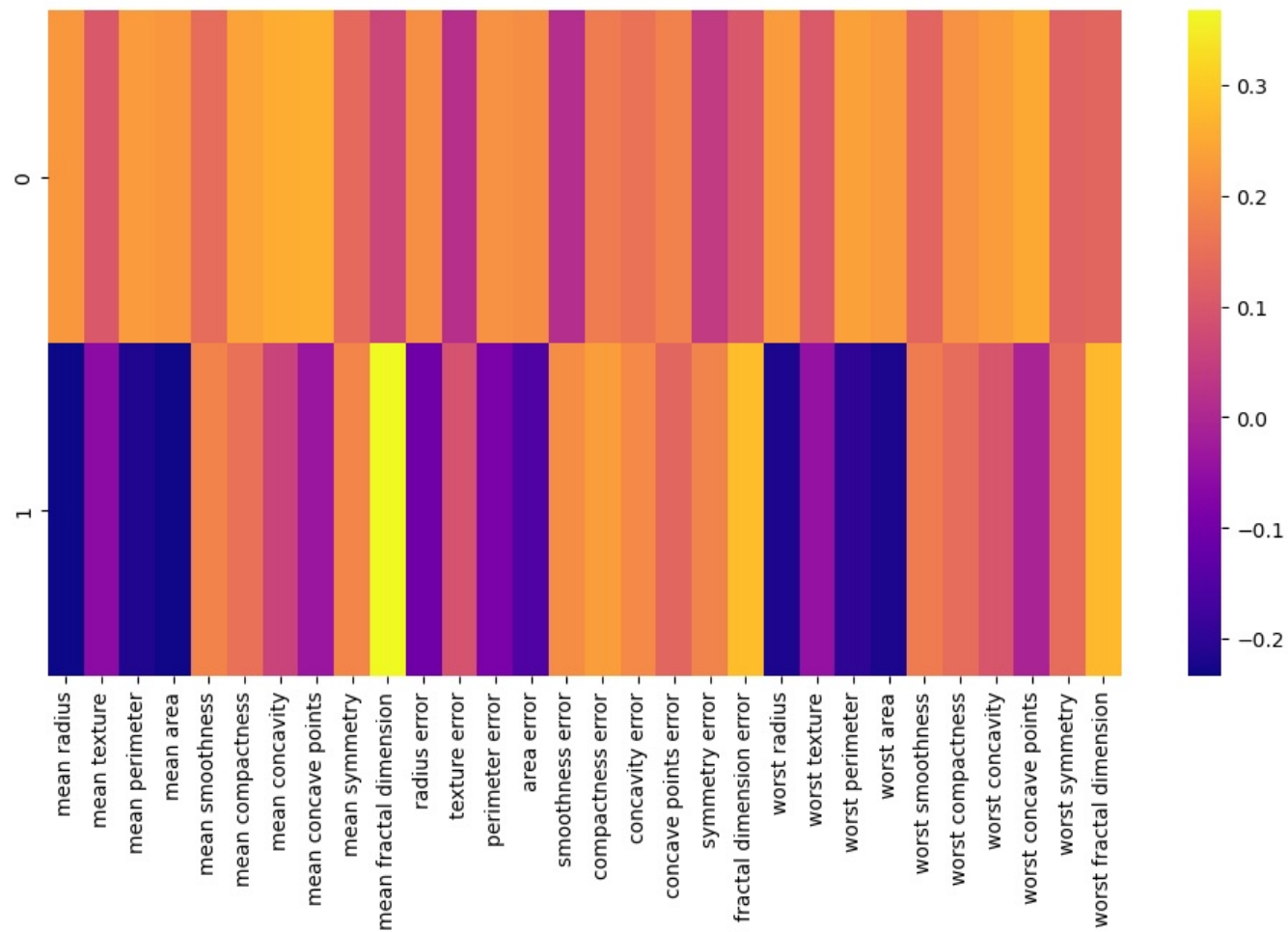


```
In [25]:  pca.components_
```

```
Out[25]:  array([[ 0.21890244,  0.10372458,  0.22753729,  0.22099499,  0.14258969,
                   0.23928535,  0.25840048,  0.26085376,  0.13816696,  0.06436335,
                   0.20597878,  0.01742803,  0.21132592,  0.20286964,  0.01453145,
                   0.17039345,  0.15358979,  0.1834174 ,  0.04249842,  0.10256832,
                   0.22799663,  0.10446933,  0.23663968,  0.22487053,  0.12795256,
                   0.21009588,  0.22876753,  0.25088597,  0.12290456,  0.13178394],
                 [-0.23385713, -0.05970609, -0.21518136, -0.23107671,  0.18611302,
                   0.15189161,  0.06016536, -0.0347675 ,  0.19034877,  0.36657547,
                  -0.10555215,  0.08997968, -0.08945723, -0.15229263,  0.20443045,
                   0.2327159 ,  0.19720728,  0.13032156,  0.183848  ,  0.28009203,
                  -0.21986638, -0.0454673 , -0.19987843, -0.21935186,  0.17230435,
                   0.14359317,  0.09796411, -0.00825724,  0.14188335,  0.27533947]])
```

```
In [27]:  df_comp = pd.DataFrame(pca.components_,columns=cancer['feature_names'])
```

```
In [29]:  plt.figure(figsize=(12,6))
          sns.heatmap(df_comp,cmap='plasma')
```

```
Out[29]:  <AxesSubplot:>
```

In [ ]: