

In [1]:

```
# maha ebrahim mohammed
# 4051350
# IA8G
# Lab 6: : Dealing with Missing Values & PCA

import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
sns.set(color_codes=True)
```

In [7]:

```
from sklearn import datasets
dataset=datasets.load_iris()

data=pd.DataFrame(dataset['data'],columns=['petal length','petal width','sepal length','sep
data['species'] = dataset['target']
data['species']=data['species'].apply(lambda x: dataset['target_names'][x])
data.head(10)
```

Out[7]:

	petal length	petal width	sepal length	sepal width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa
5	5.4	3.9	1.7	0.4	setosa
6	4.6	3.4	1.4	0.3	setosa
7	5.0	3.4	1.5	0.2	setosa
8	4.4	2.9	1.4	0.2	setosa
9	4.9	3.1	1.5	0.1	setosa

In [8]:

```
data.describe()
```

Out[8]:

	patel length	petal width	sepal length	sepal width
<b>count</b>	150.000000	150.000000	150.000000	150.000000
<b>mean</b>	5.843333	3.057333	3.758000	1.199333
<b>std</b>	0.828066	0.435866	1.765298	0.762238
<b>min</b>	4.300000	2.000000	1.000000	0.100000
<b>25%</b>	5.100000	2.800000	1.600000	0.300000
<b>50%</b>	5.800000	3.000000	4.350000	1.300000
<b>75%</b>	6.400000	3.300000	5.100000	1.800000
<b>max</b>	7.900000	4.400000	6.900000	2.500000

In [9]:

```
data.isnull().sum()
```

Out[9]:

```
patel length    0
petal width     0
sepal length    0
sepal width     0
species         0
dtype: int64
```

In [16]:

```
modData = data.append({'patel length' : np.nan, 'petal width' : 3.6, 'sepal length' : 0,
                      'sepal width' : 0.2, 'species' : 'setosa'}, ignore_index=True)
modData.describe()
```

Out[16]:

	patel length	petal width	sepal length	sepal width
<b>count</b>	150.000000	151.000000	151.000000	151.000000
<b>mean</b>	5.843333	3.060927	3.733113	1.192715
<b>std</b>	0.828066	0.436650	1.785785	0.764033
<b>min</b>	4.300000	2.000000	0.000000	0.100000
<b>25%</b>	5.100000	2.800000	1.550000	0.300000
<b>50%</b>	5.800000	3.000000	4.300000	1.300000
<b>75%</b>	6.400000	3.350000	5.100000	1.800000
<b>max</b>	7.900000	4.400000	6.900000	2.500000

In [20]:

```
print('Columns with missing values')
print(modData.isnull().sum())
print('\n columns with zero values')
print((modData[['petal length', 'petal width', 'sepal length', 'sepal width', 'species']] == 0).sum())
```

Columns with missing values

```
petal length    1
petal width     0
sepal length    0
sepal width     0
species         0
dtype: int64
```

columns with zero values

```
petal length    0
petal width     0
sepal length    1
sepal width     0
species         0
dtype: int64
```

In [21]:

```
modData[['petal length', 'petal width', 'sepal length', 'sepal width', 'species']] = modData[['petal length', 'petal width', 'sepal length', 'sepal width', 'species']]
print('Columns with missing values')
print(modData.isnull().sum())
```

Columns with missing values

```
petal length    1
petal width     0
sepal length    1
sepal width     0
species         0
dtype: int64
```

In [27]:

```
import warnings
warnings.simplefilter(action='ignore', category=FutureWarning)

modData.fillna(modData.mean(), inplace=True)
print(modData.isnull().sum())
```

```
petal length    0
petal width     0
sepal length    0
sepal width     0
species         0
dtype: int64
```

```
import warnings
warnings.simplefilter(action='ignore', category=FutureWarning)

modData.fillna(modData.median(), inplace=True)
print(modData.isnull().sum())
```

```
from sklearn.preprocessing import StandardScaler, MinMaxScaler
features = ['sepal length', 'sepal width', 'petal length', 'petal width']
x = PCA_df.loc[:, features].values
y = PCA_df.loc[:, ['target']].values
x = StandardScaler().fit_transform(x)
print(x)
```

```
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
principalComponents = pca.fit_transform(x)
principalDf = pd.DataFrame(data = principalComponents
                           , columns = ['principal component 1', 'principal component 2'])
```

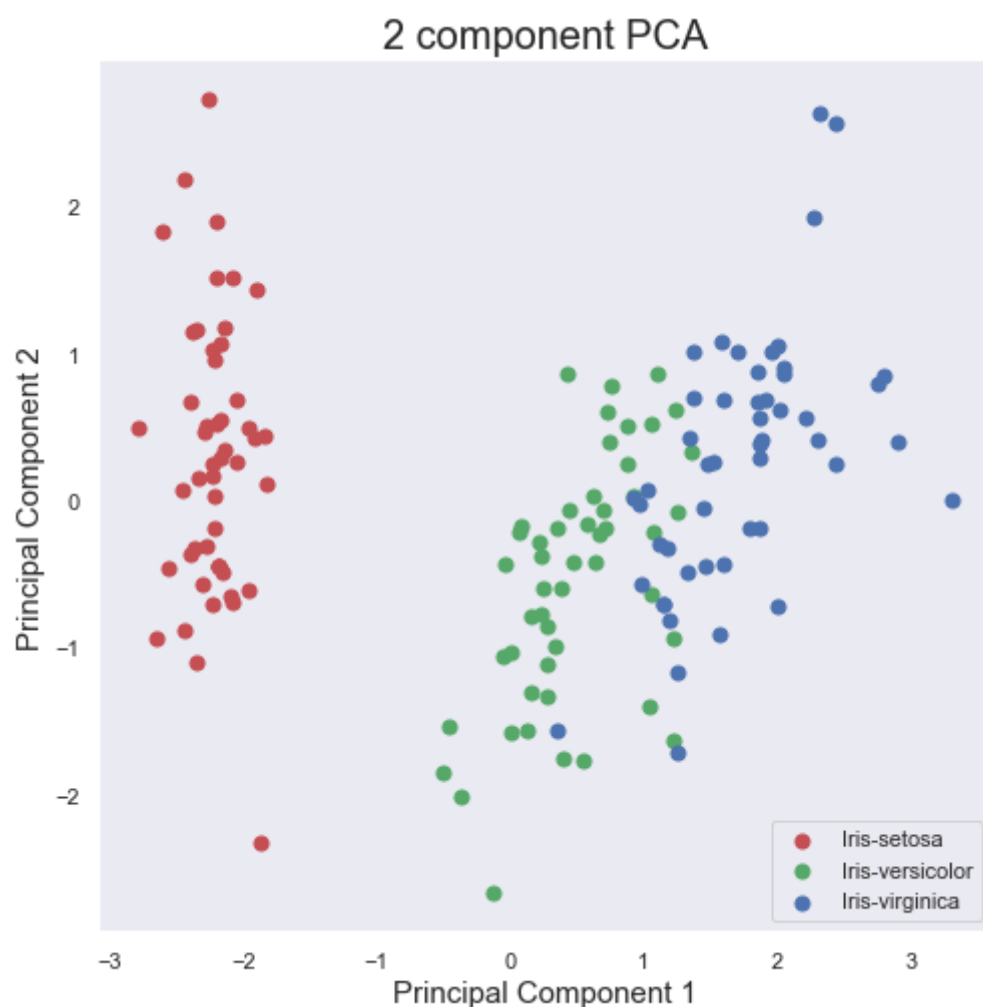
In [6]:

```
finalDf = pd.concat([principalDf, PCA_df[['target']]], axis = 1)
```

In [11]:

```
import matplotlib.pyplot as plt
%matplotlib inline
sns.set(color_codes=True)

fig = plt.figure(figsize = (8,8))
ax = fig.add_subplot(1,1,1)
ax.set_xlabel('Principal Component 1', fontsize = 15)
ax.set_ylabel('Principal Component 2', fontsize = 15)
ax.set_title('2 component PCA', fontsize = 20)
targets = ['Iris-setosa', 'Iris-versicolor', 'Iris-virginica']
colors = ['r', 'g', 'b']
for target, color in zip(targets, colors):
    indicesToKeep = finalDf['target'] == target
    ax.scatter(finalDf.loc[indicesToKeep, 'principal component 1'],
              finalDf.loc[indicesToKeep, 'principal component 2'],
              c = color
              , s = 50)
ax.legend(targets)
ax.grid()
```



In [12]:

```
print(pca.explained_variance_ratio_)
```

```
[0.72770452 0.23030523]
```

In [ ]: