

# Project By-Prince k Singh

Data Mining and Neural Networks

Breast cancer diagnosis – 2

Use Breast Cancer Wisconsin (Diagnostic) Data Set you have used in the first task and perform principal component analysis and k-means clustering.

## 1.Principle component analysis of data

first we import data set and know about data before PCA analysis Hence import required library

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
df = pd.read_csv("breast_cancer_wisconsin.csv")
df.head()
```

```
Out[1]:
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean	...	radius_worst	texture_worst	perimeter_worst	area_worst
0	842502	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	...	17.33	184.60	2019.0		
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	...	23.41	158.80	1956.0			
2	8430093	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	...	25.53	152.50	1709.0			
3	8454301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	...	26.50	98.87	567.7			
4	8435802	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	...	16.67	152.20	1575.0			

5 rows × 33 columns

As from data it is clear that it unsupervised kind of problem in which we have to diagnosis breast cancer which is MALIGANEC OR BENINE means cancerous or not cancerous so we separate this features from whole data set and also convert M AND B into 0 and 1 means numerical label with the help of label encoder

## Now it's an unsupervised learning dataset

(Because we've removed the diagnosis label) - Use this version.

```
In [2]: df_new = df.drop(['diagnosis'], 'id', 'Unnamed: 32'], axis=1)
df_new.head()
```

```
Out[2]:
```

	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean	...	radius_worst	texture_worst	perimeter_worst	area_worst
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	...	17.33	184.60	2019.0		
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	...	23.41	158.80	1956.0		
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	...	25.53	152.50	1709.0		
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	...	26.50	98.87	567.7		
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	...	16.67	152.20	1575.0		

5 rows × 30 columns

## Perform PCA on your dataset

we have need to standardize your data before PCA.

```
In [3]: from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA

#Standardize the data
scaler = StandardScaler()
norm_data = scaler.fit_transform(df_new)
df_pca = PCA(30)

#Create pca instance
# The parameter = 30 because I want my all principle components as total features is 30
pca = PCA(30)

#Fit on data
df_pca = pca.fit(norm_data)

#Access value and vectors:
eigenvectors = pca.components_
eigenvalues = pca.explained_variance_

#transform data:
#right now we have the eigenvectors and values calculated
#for the values fed in, but we haven't actually generated
#principal component variables. Here is when we create them:
pca_vals = pca.transform(norm_data)
pca_vals
df_pca

Out[5]: PCA(n_components=30)
```

Hence I present all eigenvalues in the form of list

```
In [6]: list(eigenvalues)
```

```
Out[6]: [13.384998794374538,
5.781374607281466,
2.822833508682246,
1.8841275177382625,
1.65133243391185,
1.269482398829874,
0.8746888817899506,
0.4774562546895078,
0.4176287821878163,
0.3513187488173296,
0.29443315349116456,
0.2612215158652074,
0.24178424323831318,
0.15728614921592923,
0.09430869560185278,
0.088083484471737674,
0.0890831305151,
0.85271142221814774,
0.8495947632229815,
0.8321242695586633,
0.82748771133894375,
0.8389256308424315,
0.82438369135459111,
0.8106076399410535,
0.815598527134418559,
0.80815287321767276,
0.808012612579184481,
0.8015921360611976581,
0.8007501214127186319,
0.8061332798566399843]
```

NOW I have to plot eigenvalue as abscissa of its number as indirectly EIGENVALUE(lambda) also represent principle components of co-variance matrix so indirectly I have to plot principle components as a function number

```
In [7]: # PCA for dimensionality reduction (non-visualization)
pca_n_components = 30
pca_data = pca.fit_transform(norm_data)

percentage_var_explained = pca.explained_variance_ / np.sum(pca.explained_variance_)
cum_var_explained = np.cumsum(percentage_var_explained)

# Plot the PCA spectrum
plt.figure(figsize=(6, 4))

plt.clf()
plt.plot(cum_var_explained, linewidth=2)
plt.axis('tight')
plt.grid()
plt.xlabel('n_components')
plt.ylabel('Cumulative explained variance')
plt.show()
```



AS in above plot it clear from KAISER LAW for retaining 90% of information we have to take principle components=6 nearly

## 2.Data visualization using principle components

now present the histogram of first 3 components

```
In [8]: #Create pca instance
# The parameter = 3 because I want my first 3 components principle components as total features is 30
pca = PCA(3)

#Fit on data
df_pca = pca.fit(norm_data)
pca_data = pca.fit_transform(norm_data)

#Access value and vectors:
eigenvectors = pca.components_
eigenvalues = pca.explained_variance_

#transform data:
#right now we have the eigenvectors and values calculated
#for the values fed in, but we haven't actually generated
#principal component variables. Here is when we create them:
pca_vals = pca.transform(norm_data)
pca_vals
df_pca

Out[8]: PCA(n_components=3)
```

now with help of encoder I convert M and B into 1 and 0

```
In [9]: def encoder(data):
    if data == 'M':
        return 1
    else:
        return 0
df_target = df['diagnosis'].apply(encoder)
df_target
```

```
Out[9]:
```

0	1
1	1
2	1
3	1
4	1
...	...
564	1
565	1
566	1
567	0
568	0

Name: diagnosis, Length: 569, dtype: int64

```
In [10]: pca_df = pd.DataFrame(pca_data, columns=['pca0', 'pca1', 'pca2'])
pca_df['diagnosis'] = df['diagnosis']
print('Shape of PCA dataset is {}'.format(pca_df.shape))
pca_df.head()
```

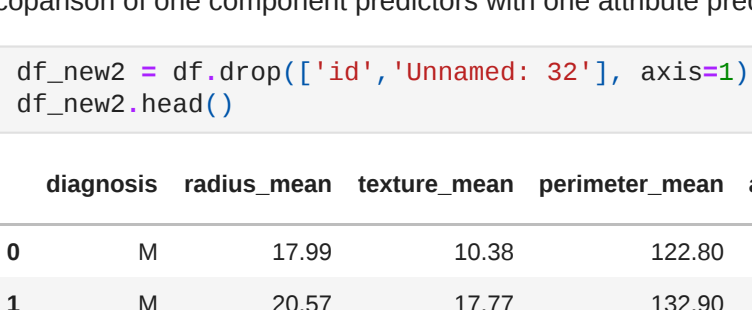
```
Out[10]:
```

	pca0	pca1	pca2	diagnosis
0	9.192837	1.948583	-1.221564	M
1	2.387802	3.768172	-0.528292	M
2	5.733906	-1.075174	0.951748	M
3	7.122953	10.275589	-0.232769	M
4	3.935302	-1.948071	1.389768	M

NOW we plot HISTOGRAM plot for 1st 3 principle component

```
In [11]: import seaborn as sns
sns.histplot(pca_df,)
```

```
Out[11]: <AxesSubplot:ylabel='Count'>
```



but from above plot it is not clear that which one separable hence we have to plot pairwise scatter plot for clear view

```
In [12]: sns.pairplot(pca_df, hue='diagnosis', markers=["o", "s"], corner=False)

Out[12]: <matplotlib.axes._subplots.PairGrid at 0xf7b48c5489>
```



from above plot it seems pca0 and pca1 are have good separation among them

coparsion of one component predictors with one attribute predictors by plotting pair plots of general attribute

```
In [13]: df_new2 = df.drop(['id', 'Unnamed: 32'], axis=1)
df_new2.head()
```

	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	...	radius_worst	texture_worst	perimeter_worst	area_worst
0	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	...	25.38	184.60	2019.0	
1	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	...	24.99	158.80	1956.0	
2	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	...	23.57	152.50	1709.0	
3	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	...	26.51	98.87	567.7	
4	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	...	14.94	152.20	1575.0	

5 rows × 31 columns

```
In [14]: #sns.pairplot(df_new2, hue='diagnosis', markers=["o", "s"], corner=False)
```

## 3.K-means clustering

now performing K-MEANS clustering for k=2,3 and 5

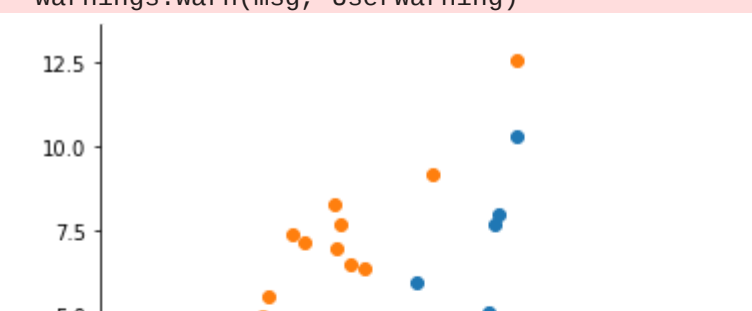
```
In [15]: from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=2)
kmeans.fit(norm_data)
y_kmeans = kmeans.predict(norm_data)
plt.scatter(norm_data[:, 0], norm_data[:, 1], c=y_kmeans, s=50, cmap='viridis')
centers = kmeans.cluster_centers_
plt.scatter(centers[:, 0], centers[:, 1], c='black', s=200, alpha=0.5)

Out[15]: <matplotlib.collections.PathCollection at 0xf7b41ec3438>
```



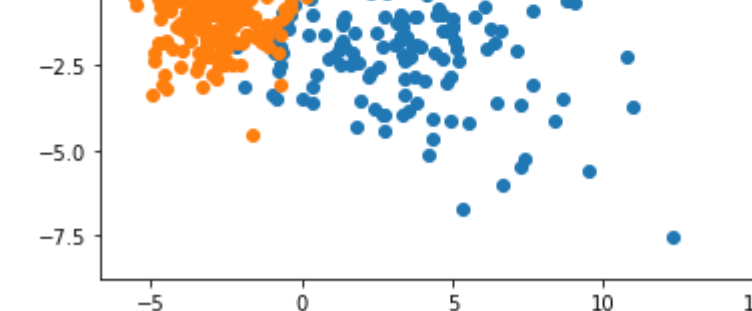
```
In [16]: kmeans = KMeans(n_clusters=3)
kmeans.fit(norm_data)
y_kmeans = kmeans.predict(norm_data)
plt.scatter(norm_data[:, 0], norm_data[:, 1], c=y_kmeans, s=50, cmap='viridis')
centers = kmeans.cluster_centers_
plt.scatter(centers[:, 0], centers[:, 1], c='black', s=200, alpha=0.5)

Out[16]: <matplotlib.collections.PathCollection at 0xf7b411333ab>
```



```
In [17]: kmeans = KMeans(n_clusters=5)
kmeans.fit(norm_data)
y_kmeans = kmeans.predict(norm_data)
plt.scatter(norm_data[:, 0], norm_data[:, 1], c=y_kmeans, s=50, cmap='viridis')
labels = kmeans.labels_
plt.scatter(centers[:, 0], centers[:, 1], c='black', s=200, alpha=0.5)

Out[17]: <matplotlib.collections.PathCollection at 0xf7b4118ab58>
```



now I am finding the davisbouldin index with centroids of clusters so first define a function of davisbouldin fifth it parameter of satandardised data predicted label and centroids hence we find individual index and centroids table as well

now I plot first two principle components and plot it with different colour I used facegrid plot for representation

```
In [18]: pca = PCA(2)

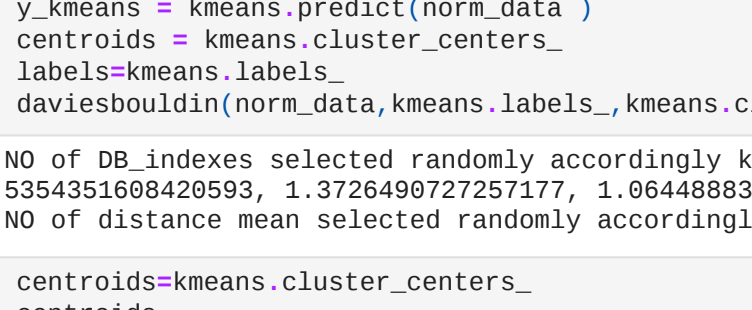
#Fit on data
df_pca = pca.fit(norm_data)
pca_data = pca.fit_transform(norm_data)
pca_df = pd.DataFrame(pca_data, columns=['1st_principal', '2nd_principal'])
pca_df['diagnosis'] = df['diagnosis']
print('Shape of PCA dataset is {}'.format(pca_df.shape))
pca_df.head()
```

```
Out[18]:
```

	1st_principal	2nd_principal	diagnosis
0	9.192837	1.948583	M
1	2.387802	-3.768172	M
2	5.733906	-1.075174	M
3	7.122953	10.275589	M
4	3.935302	-1.948072	M

```
In [19]: import seaborn as sns
sns.scatterplot(pca_df, hue='diagnosis', size=60).map(plt.scatter, '1st_principal', '2nd_principal').add_legend()
plt.show()
```

C:\ProgramData\Anaconda3\Lib\site-packages\seaborn\axisgrid.py:316: UserWarning: The 'size' parameter has been renamed to 'height'; please update your code. warnings.warn(msg, UserWarning)



## daviesbouldin index

we calculate daviesbouldin index and first define function and get out put of davies index list and distance mean

```
In [22]: def daviesbouldin(X, labels, centroids):
    from scipy.spatial.distance import pdist, euclidean

    nbre_of_clusters = len(centroids) #Get the number of clusters
    distances = []
    for i, e in enumerate(X): #Store the (nbre_of_clusters) distance intra-cluster distances by cluster
        distances.append(euclidean(X[e], centroids[labels[e]]))
    distances_means = [] #Store the mean of these distances
    DB_indexes = [] #Store Davies-Bouldin index of each pair of cluster
    second_cluster_idx = 0 #Store index of the second cluster of each pair
    first_cluster_idx = 0 #Store index of first cluster of each pair to 0

    # Step 1: Compute euclidean distances between each point of a cluster to their centroid
    for cluster in range(nbre_of_clusters):
        for point in range(nbre_of_clusters):
            distances[cluster].append(euclidean(X[labels == cluster][point], centroids[cluster]))

    # Step 2: Compute the mean of these distances
    for e in distances:
        distances_means.append(np.mean(e))

    # Step 3: Compute euclidean distances between each pair of centroid
    ctrds_distances = pdist(centroids)

    # Step 4: Compute Davies-Bouldin index of each pair of cluster
    for i, e in enumerate(X): #Store the (nbre_of_clusters) distance intra-cluster distances by cluster
        second_cluster_idx.append(e)
        if second_cluster_idx[i] != nbre_of_clusters - 1:
            first_cluster_idx.append(i)
            DB_indexes.append((distances_means[first_cluster_idx[i]] + distances_means[e]) / ctrds_distances[i])

    # Step 5: Compute the mean of all DB_indexes
    print("NO of DB_indexes selected randomly according k-means:", DB_indexes)

    print("NO of distance mean selected randomly according k-means:", distances_means)
```

```
In [23]: kmeans = KMeans(n_clusters=5)
kmeans.fit(norm_data)
y_kmeans = kmeans.predict(norm_data)
centroids = kmeans.cluster_centers_
labels = kmeans.labels_
daviesbouldin(norm_data, kmeans.labels_, kmeans.cluster_centers_)
```

NO of DB\_indexes selected randomly according k-means: [0.8271299466951883, 0.893719757932939, 1.17431535862191681, 1.478486562084223, 0.7320684422841057, 1.5351351698420583, 1.372648972757177, 1.06448885683282, 1.0180611063617975, 1.4874870729275883]

NO of distance mean selected randomly according k-means: [5.31678131343587, 2.8670397471503184, 3.4492761556104274, 5.2315312146891376]

```
In [24]: centroids=kmeans.cluster_centers_

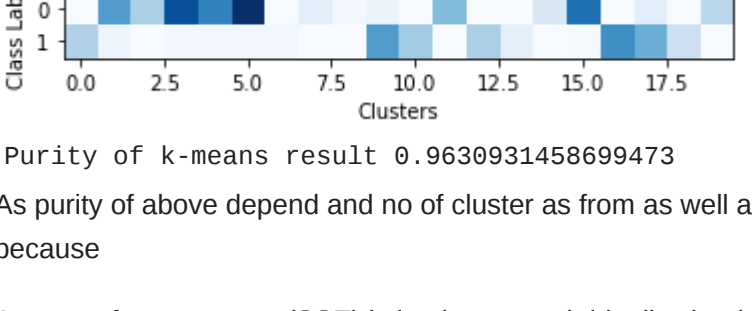
Out[24]: array([[ -7.12937056e-01, -3.41133570e-01, -6.96841381e-01,
-6.83437556e-01, 3.75621627e-01, -1.56232478e-01,
-1.84646074e-01, 4.35731244e-01, 9.39518151e-02,
4.58258603e-01, -3.97749856e-01, 1.08897298e-01,
-3.79158225e-01, -4.24342764e-01, 5.45371508e-01,
-5.18416167e-02, -1.25841235e-01, -7.88511414e-02,
1.73267523e-01, 0.10193188e-01, -7.03178268e-01,
-3.68641230e-01, -6.84743232e-01, -6.40978285e-01,
3.62343386e-01, -2.68487735e-01, -3.80691954e-01,
-2.68952787e-01, -1.34845389e-01, 6.84106730e-02,
2.63236185e-01, 7.98578088e-01, 2.11878771e+00,
2.21216632e+00, 8.98971322e-01, 1.97017715e+00,
2.16779361e+00, 2.23855288e+00, 3.98728354e+00,
4.40485733e-01, 2.43894161e+00, 2.69884175e-01,
2.54079204e+00, 2.40878091e+00, 2.80727271e-01,
1.22718599e+00, 9.71988922e-01, 1.27879725e+00,
5.06544207e-01, 6.33740260e-01, 2.11722858e+00,
5.89462639e-01, 2.22382738e+00, 2.26515289e+00,
3.74797240e-01, -9.99891320e-01, -6.88076950e-01,
-7.14545387e-01, -4.58772071e-01, -6.42392777e-01,
1.35949171e+00, 4.66614847e-01, 1.12548043e+00,
1.63487862e+00, -5.54073590e-01, 1.24190506e-01,
6.57846267e-01, 9.08773862e-01, 1.15191588e-01,
-5.72212720e-01, 6.47852224e-01, 9.37209356e-02,
5.82957545e-01, 6.38593166e-01, -2.5969275e-01,
2.25431977e-02, 1.29424452e-01, 4.64589744e-01,
-1.98253107e-01, -1.74387454e-01, 1.38426950e+00,
4.96518972e-01, 2.16839888e+00, 1.10281971e+00,
2.39383324e-01, 1.68449141e-01, 1.79915250e+00,
9.12583189e-01, 2.74809277e-01, -8.07855787e-02,
-2.08926323e-01, 2.56895250e-01, -1.12977573e-01,
-2.48373266e-01, 9.97228898e-01, 1.38458236e+00,
1.33891527e+00, 6.45294977e-01, 1.93851208e+00,
1.63487862e+00, -5.54073590e-01, 1.24190506e-01,
7.63584952e-02, -1.41973688e-01, 5.12832420e-01,
1.40683322e+00, 1.31329398e+00, 9.79254295e-01,
5.07281340e-01, 1.49675684e+00, -1.03934301e-01,
4.08458350e-01, -1.12597515e-03, -1.70797024e-01,
1.12512289e+00, 1.16581750e+00, 1.47788930e+00,
9.25144181e-01, 1.06481117e+00, 1.88221983e+00]])
```

## 4. Clustering and Classification

now I find purty which is kind of accuracy in kmeans 1st of all plot sum of squared error with and it seems after 15 cluster the rate of decrease of error becomes constant hence we take no of cluster 15 and then find purty as it is random hence each time purty can be range as well

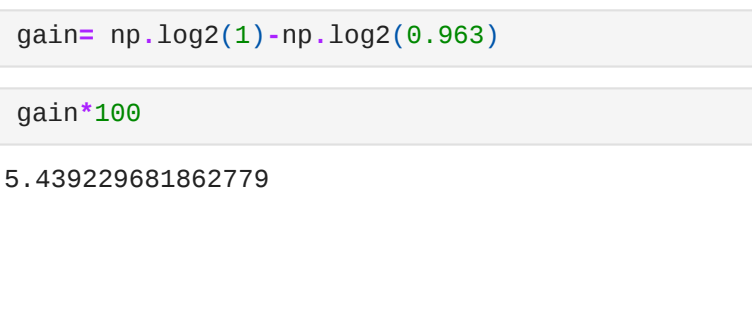
```
In [25]: sse = []
for k in range(1, 30):
    kmeans = KMeans(n_clusters=k, max_iter=1000).fit(df_new)
    #print(data['cluster'])
    sse[k] = kmeans.inertia_ # Inertia: Sum of distances of samples to their closest cluster center
    plt.plot(list(sse.keys()), list(sse.values()))
    plt.xlabel('number of cluster')
    plt.ylabel('sse')
    plt.show()

C:\Users\Himov\AppData\Roaming\Python\Python38\site-packages\sklearn\cluster\_kmeans.py:881: UserWarning: KMeans is known to have a memory leak on Windows w
th MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=3.
warnings.warn(
```



```
In [30]: kmeans = KMeans(n_clusters=29)
kmeans.fit(norm_data)
y_kmeans = kmeans.predict(norm_data)
def evaluate(y_true, y_clusters):
    cm = np.zeros((np.max(y_true)+1, np.max(y_clusters)+1))
    for i in range(y_true.size):
        only_true[i], y_clusters[i] += 1
    plt.imshow(cm, interpolation='none', cmap=plt.cm.Blues)
    plt.ylabel('Class Labels')
    plt.xlabel('Clusters')
    plt.show()
    purity = 0.
    for cluster in range(cm.shape[1]): # clusters are along columns
        purity += np.max(cm[:, cluster])
    return purity/y_true.size

purity = evaluate(df_target, y_kmeans)
print('Purity of k-means result', purity)
```



Purity of k-means result 0.963093458699473

As purity of above depend on no of cluster as from as well as it is also change with each run and because of randomness in clustering process it should be ideally 1 but in real practice it is closer to 1

As sum of square error (SSE) in above graph ideally should be 0 or touches 0 but in mathematical it takes infinite iteration so we loss of information from target attribute over categorical attribute or in other word we relatively gain information from categorical to target

## relative information gain

```
information gain = information(target) - information(categoryal)
relative information gain = log(purity value target) - log(purity value categorical)
here in above case purity value target = 1, purity value categorical = 0.963

In [41]: gain = np.log2(1) - np.log2(0.963)

Out[43]: gain=186

In [43]: 5.439229681862779
```