

Q-1 : PCA

- Akshay Bankar (2019201011)

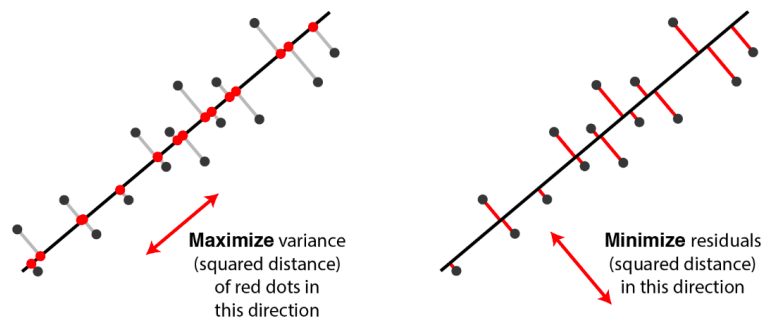
PCA as dimensionality reduction technique:

To reduce the data dimensionality we should be able to summarize the data with less characteristics (features). PCA does this by taking the linear combination of the existing features and constructs some new features that are good alternative representation of the original data.

PCA tries to pick the best direction, referred to as principal components and projects the points onto these directions.

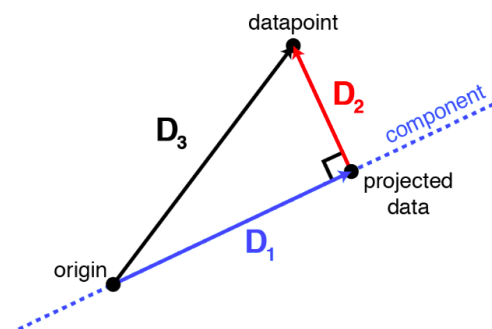
How does PCA do this :

There are two ways to formulate such projection. First, look for some features that strongly differ across data points, thus, PCA looks for features that captures as much variation across data points as possible. Second, that we are looking for the features that would allow us to "reconstruct" the original features.



Two equivalent views of principal component analysis.

From the following figure it can be seen that these two objectives are same. It is equivalent to either maximize remaining variance or minimize lost variance to find the principal components.



$$D_3^2 = D_1^2 + D_2^2$$

initial variance = remaining variance + lost variance

$$\|a_i\|^2 = \|w_i c\|^2 + \|a_i - w_i c\|^2$$

this is constant maximize this or minimize this

```
In [1]: 1 import numpy as np
        2 from sklearn.decomposition import PCA
        3 import cv2
        4 import glob
```

PCA formulation :

given X , the centered data matrix, the projection, Xw (dot product between the data point, X and the projection weight, w), its variance can be computed as follows:

$$\frac{1}{n-1} (Xw)^T Xw = w^T \left(\frac{1}{n-1} X^T X \right) w = w^T C w$$

Where C is the covariance matrix of the data X . The covariance matrix is a $d \times d$ matrix (d is the total number of dimensions, features) where each element represents the covariance between two features. The covariance between two features x_j and x_k is calculated as follows:

$$\sigma_{jk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j) (x_{ik} - \bar{x}_k)$$

$$\bar{x}_j$$

is the mean of vector (feature)

$$\bar{x}_j = \sum_{i=1}^n x_{ij}$$

The covariance matrix is then,

$$C = \frac{1}{n-1} X^T X$$

Now the objective function for PCA is

$$\begin{aligned} &\underset{w}{\text{maximize}} && w^T C w \\ &\text{subject to} && w^T w = 1 \end{aligned}$$

This objective function can be solved by the Lagrange multiplier, minimizing the loss function:

$$\begin{aligned} \frac{\partial L}{\partial w} &= Cw - \lambda w = 0 \\ &= Cw = \lambda w \end{aligned}$$

After solving the equation above, we'll obtain eigenvector and eigenvalue pairs, where every eigenvector has a corresponding eigenvalue. An eigenvector is essentially the direction of each principal component and the eigenvalue is a number telling us how much variance there is in the data in that direction, in other words, how spread out the data is on the line.

Load the dataset

In [2]: `img_files = glob.glob('dataset/*')`

Read the color image file.

Resize the image to size (64,64)

Convert image to gray scale

Flatten the data to obtain row vectors

```
In [82]: 1 gray_images = []
2 labels = []
3 for file in img_files:
4     img = cv2.imread(file)
5     img = cv2.resize(img,(64,64),interpolation=cv2.INTER_AREA) #None, fx=0.5
6     flat_img = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY).flatten()
7     gray_images.append(flat_img)
```

Standardize the data

```
In [84]: 1 n_features = np.shape(gray_images[0])[0]
2 #print(n_features)
3 std_dev = np.zeros(np.shape(gray_images))
4 print(std_dev.shape)
5 gray_images = np.asarray(gray_images)
6 for i in range(n_features):
7     std_dev[:,i] = (gray_images[:,i] - np.mean(gray_images[:,i]))/gray_images[:,i]
```

(520, 4096)

Obtain the covariance matrix

```
In [85]: 1 #print(np.shape(gray_images))
2 # Covariance
3 #np.set_printoptions(precision=3)
4 cov = np.cov(std_dev.T)# Eigen Values
```

Obtain the Eigenvectors and Eigenvalues from the covariance matrix

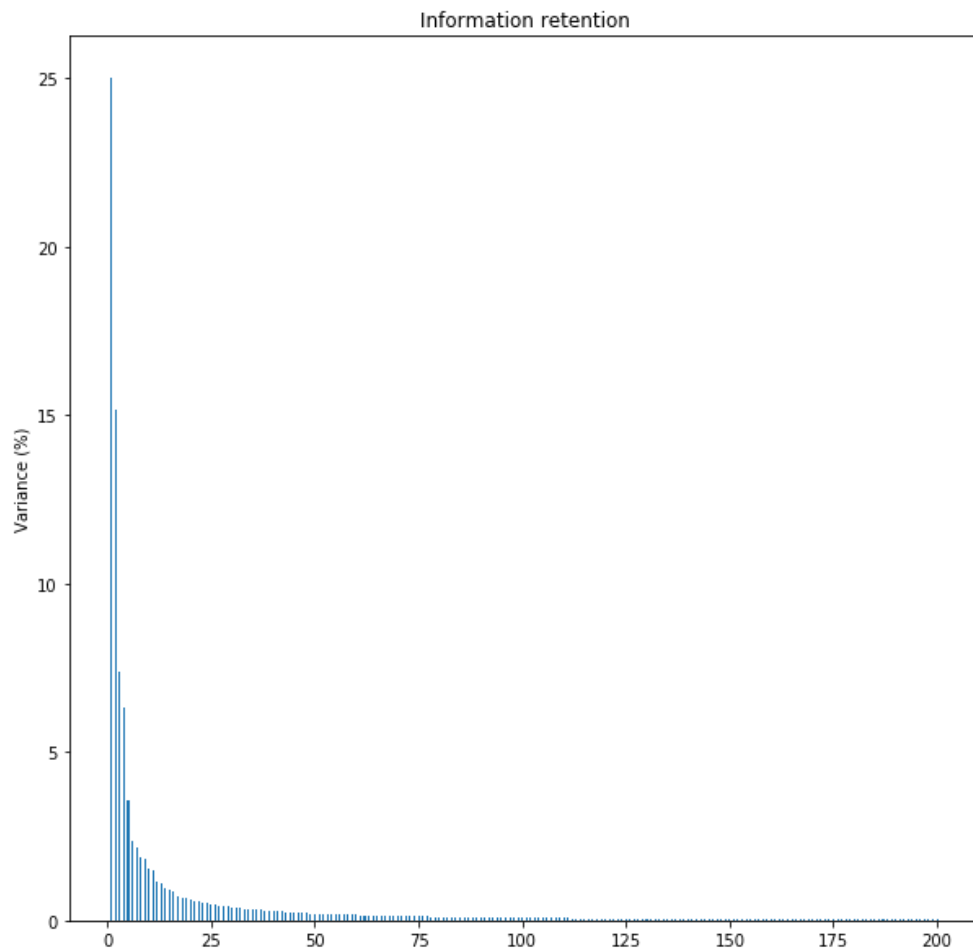
```
In [86]: 1 from scipy.linalg import eigh
2 EigVal,EigVec = eigh(cov)
3 #EigVal,EigVec = np.linalg.eig(cov)
4 #print("Eigenvalues:\n\n" EigVal "\n")
```

Sort eigenvalues and corresponding eigenvectors in descending order and compute the principal components

```
In [87]: 1 order = EigVal.argsort()[::-1]
2 EigVal = EigVal[order]
3 EigVec = EigVec[:,order] #Projecting data on Eigen vector directions results
4 PC = np.matmul(std_dev, EigVec) #cross product
```

Plot of graph which shows the amount of information contained in each direction of eigenvectors

```
In [88]: 1 import matplotlib.pyplot as plt
2 import matplotlib.gridspec as grid
3
4 plt.figure(figsize=(10,10))
5 plt.bar([i+1 for i in range(200)],EigVal[:200]/sum(EigVal)*100,align='center')
6 plt.ylabel('Variance (%)')
7 plt.title('Information retention');
```



Number of components that contain 80% of information

```
In [89]: 1 percent=0
2 count=0
3 thresh = 0.8
4 for i in EigVal:
5     percent+= i/sum(EigVal)
6     count+=1
7     if percent > thresh:
8         break;
9 print("Number of components that contain 80% of information : " count)
```

Number of components that contain 80% of information : 27

Reconstruction image using these components containing 80% of information

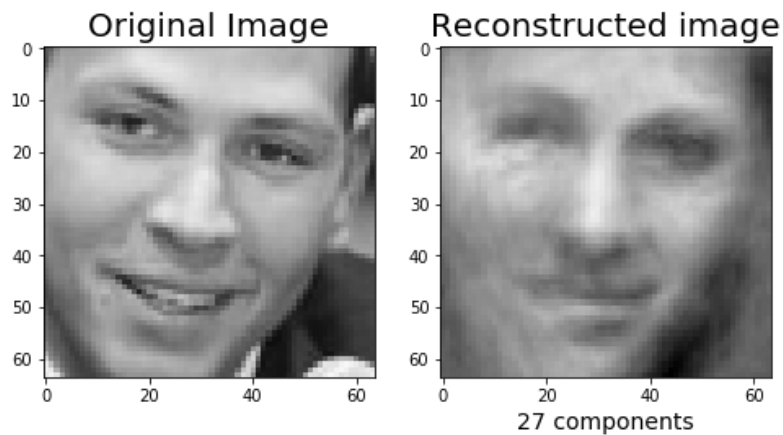
```

In [90]: 1 import matplotlib.pyplot as plt
          2
          3 plt.figure(figsize=(8,4));
          4
          5 # Original Image
          6 plt.subplot(1, 2, 1);
          7 plt.imshow(gray_images[0].reshape(64,64),
          8           cmap = plt.cm.gray, interpolation='nearest',
          9           clim=(0, 255));
         10 #plt.xlabel('64x64 components', fontsize = 14)
         11 plt.title('Original Image', fontsize = 20);
         12
         13 PC_img = np.matmul(gray_images, EigVec[:, :count])
         14 print(PC_img.shape)
         15 reconstructed_img = np.matmul(PC_img, EigVec[:, :count].T)
         16 plt.subplot(1, 2, 2);
         17 plt.imshow(reconstructed_img[0].reshape(64, 64),
         18           cmap = plt.cm.gray, interpolation='nearest',
         19           clim=(0, 255));
         20 plt.title('Reconstructed image', fontsize = 20);
         21 plt.xlabel('27 components', fontsize = 14)
         22

```

(520, 27)

Out[90]: Text(0.5, 0, '27 components')



Plot to show total mean square error over all train images vs the number of principal components

```

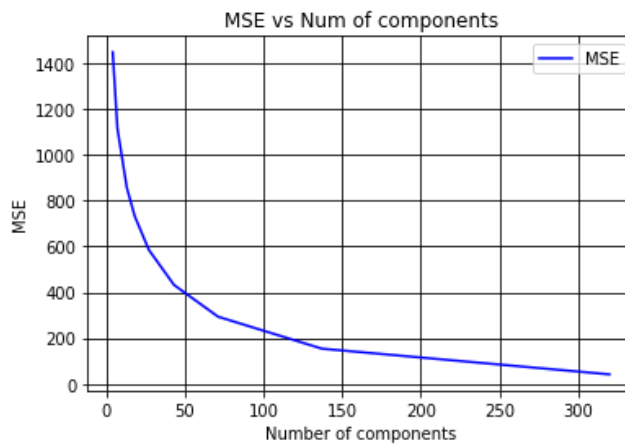
In [91]: 1 from sklearn.metrics import mean_squared_error
2
3 pca_comp = [0.99, 0.95, 0.90, 0.85, 0.80, 0.75, 0.70, 0.60, 0.50]
4 n_comp = []
5 mse = []
6 for num in pca_comp:
7     count = 0
8     percent = 0
9     for i in EigVal:
10         percent+= i/sum(EigVal)
11         count+=1
12         if percent > num:
13             break;
14     n_comp.append(count)
15     print("Number of components with", num, "% of information :",count)
16     PC_img = np.matmul(gray_images, EigVec[:, :count])
17     #print(PC_img.shape)
18     reconstructed_img = np.matmul(PC_img, EigVec[:, :count].T)
19     mse.append(mean_squared_error(gray_images, reconstructed_img))
20
21 plt.title("MSE vs Num of components")
22 plt.ylabel('MSE')
23 plt.xlabel('Number of components')
24
25 plt.plot(n_comp,mse,'b', label="MSE")
26 #plt.plot(k,manhatt_success,'g', label="Manhatt. succ. rate")
27 plt.legend()
28 plt.grid(True, color='k')
29 plt.show()

```

```

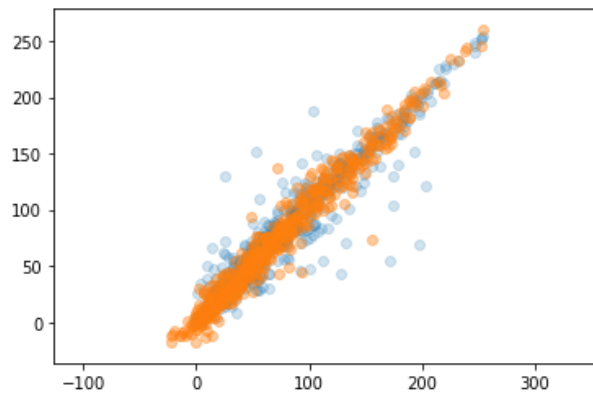
Number of components with 0.99 % of information : 320
Number of components with 0.95 % of information : 137
Number of components with 0.9 % of information : 71
Number of components with 0.85 % of information : 43
Number of components with 0.8 % of information : 27
Number of components with 0.75 % of information : 18
Number of components with 0.7 % of information : 13
Number of components with 0.6 % of information : 7
Number of components with 0.5 % of information : 4

```

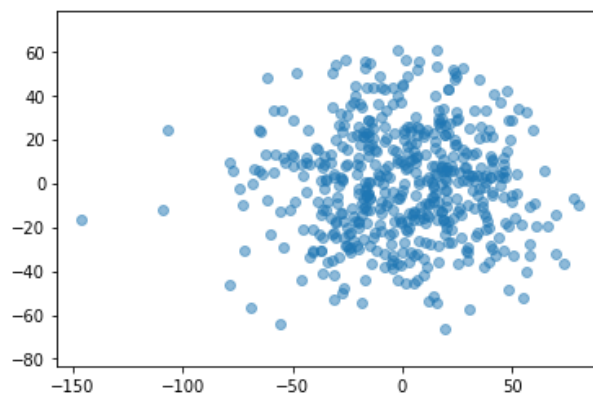


1D, 2D and 3D plots of the principal components

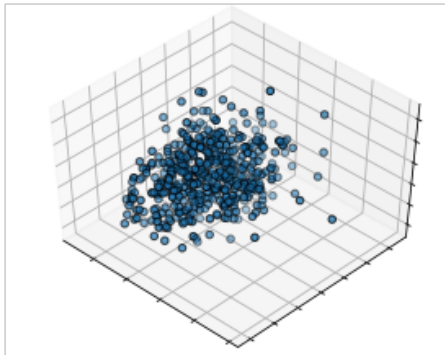
```
In [92]: 1 PC_img = np.matmul(gray_images, EigVec[:, :137])
2 reconstructed_img = np.matmul(PC_img, EigVec[:, :137].T)
3 plt.scatter(gray_images[:, 0], gray_images[:, 1], alpha=0.2)
4 plt.scatter(reconstructed_img[:, 0], reconstructed_img[:, 1], alpha=0.4)
5 plt.axis('equal').
```



```
In [93]: 1 def draw_vector(v0, v1, ax=None):
2     ax = ax or plt.gca()
3     arrowprops=dict(arrowstyle='->',
4                     linewidth=2,
5                     shrinkA=0, shrinkB=0)
6     ax.annotate('', v1, v0, arrowprops=arrowprops)
7
8     # plot data
9     plt.scatter(PC[:, 0], PC[:, 1], alpha = 0.5)
10
11 for length, vector in zip(EigVal[0:2], EigVec[:, 0:2]):
12     v = vector * 3 * np.sqrt(length)
13     #print(np.shape(np.mean(gray_images[:, :2], axis=0)))
14     draw_vector(np.mean(gray_images[:, :2], axis=0), np.mean(gray_images[:, :2]
15 plt.axis('equal');
```



```
In [94]: 1 import matplotlib.pyplot as plt
2 from mpl_toolkits.mplot3d import Axes3D
3
4 fig = plt.figure(1, figsize=(4, 3))
5 plt.clf()
6 ax = Axes3D(fig, rect=[0, 0, .95, 1], elev=48, azimuth=134)
7
8 plt.cla()
9 ax.scatter(PC[:, 0], PC[:, 1], PC[:, 2], cmap=plt.cm.nipy_spectral,
10           edgecolor='k')
11
12 ax.w_xaxis.set_ticklabels([])
13 ax.w_yaxis.set_ticklabels([])
14 ax.w_zaxis.set_ticklabels([])
15
16 plt.show()
```



In []: 1