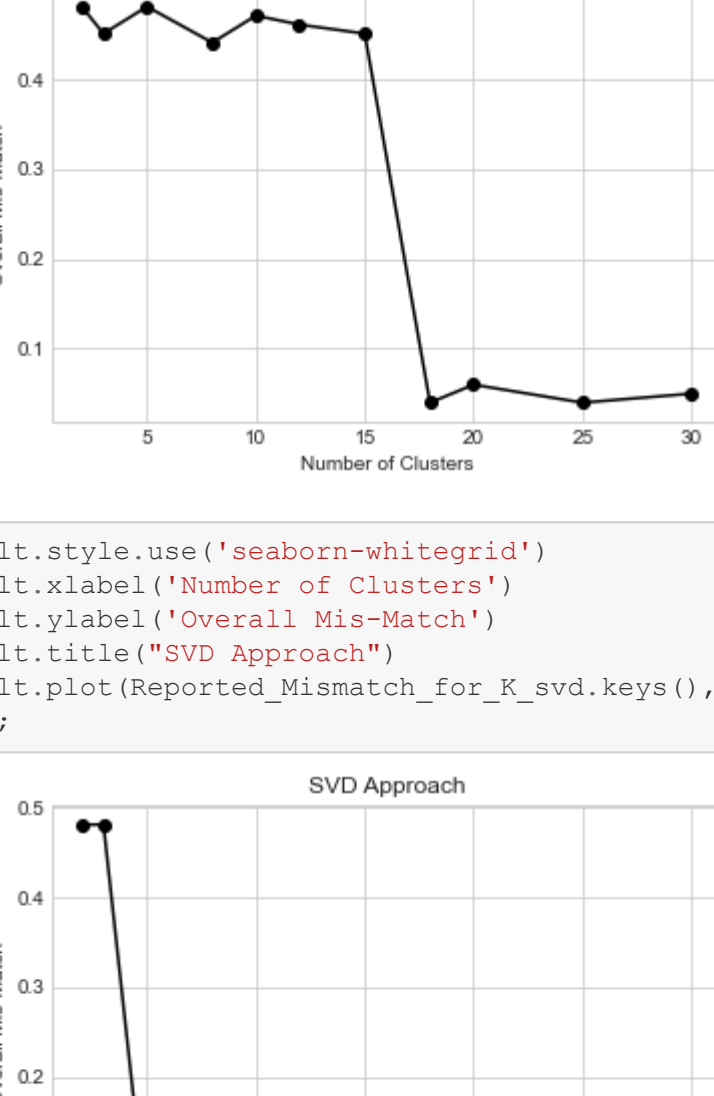
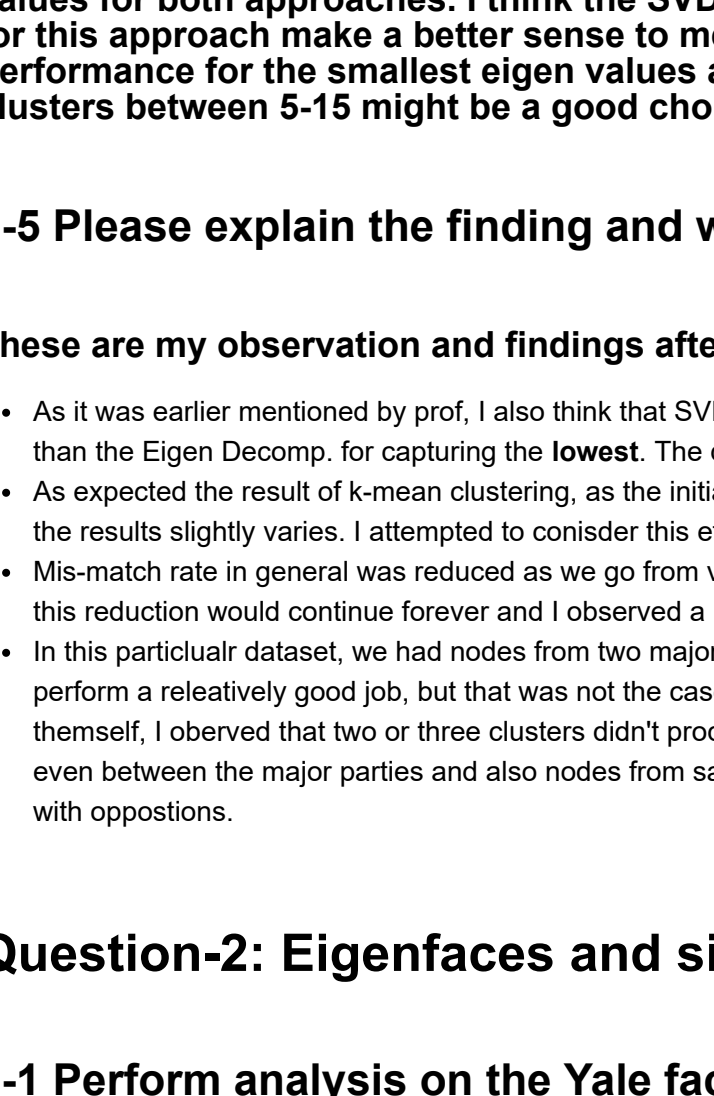



```
[41]: import matplotlib.pyplot as plt
import matplotlib.pyplot as plt
plt.style.use('seaborn-whitegrid')
plt.xlabel('Number of Clusters')
plt.ylabel('Overall Mis-Match')
plt.title('Eigen Decomp. Approach')
plt.plot(Reported_Mismatch_for_K_keys(), Reported_Mismatch_for_K_values(), '-ok', color='black')
In [42]:
```



```
In [42]: plt.style.use('seaborn-whitegrid')
plt.xlabel('Number of Clusters')
plt.ylabel('Overall Mis-Match')
plt.title('SVD Approach')
plt.plot(Reported_Mismatch_for_K_keys(), Reported_Mismatch_for_K_values(), '-ok', color='black')
In [43]:
```



I have tested Eigen Decomp. and SVD approaches and plotted the mis-matches in respect with K values for both approaches. I think the SVD approach produced more trustworthy results. The plot for this approach make a better sense to me, also there are reported issues for the Eigen Decomp. performance for the smallest eigen values and vectors. I hence based on the results I think number of clusters between 5-15 might be a good choice.

1-5 Please explain the finding and what can you learn from this data analysis.

These are my observation and findings after doing question-1:

- As it was earlier mentioned by prof. I also think that SVD is more reliable to derive sorted eigen values and capturing the highest values than the Eigen Decomp. for capturing the lowest. The derived mis-match via these two approaches suggest that.
- As expected the result of k-mean clustering, as the initial points were randomly selected, are not unique and hence after each repeat the results slightly varies I attempted to consider this effect and reported sort of averaged statistics.
- Mis-match rate in general was reduced as we go from very low number of clusters to higher numbers, but I noticed that didn't mean that this reduction would continue forever and I observed a reverse effect for cluster numbers more than 20.
- In this particular dataset, we had nodes from two major political groups, maybe then the first exception was the two clusters may perform a relatively good job, but that was not the case, even when I removed the two nodes that were only in communication with themselves, I observed that two or three clusters didn't produce a satisfactory results. That suggests there were considerable connections even between the major parties and also nodes from same parties may have some political disagreements, or considerable connection with oppositions.

Question-2: Eigenfaces and simple face recognition

2-1 Perform analysis on the Yale face dataset for Subject 1 and Subject 2, respectively, using all the images EXCEPT for the two pictures named subject01-test.gif and subject02-test.gif. Plot the first 6 eigenfaces for each subject. When visualizing, please reshape the eigenfactors into proper images. Please explain can you see any patterns in the top 6 eigenfaces?

```
In [43]: from PIL import Image
import numpy as np
import matplotlib.pyplot as plt
from os import listdir
from os.path import isfile, join
import math
import scipy
```

```
In [44]: def read_image(file_path):
        return Image.open(file_path)

def downsample_image(image):
    w, h = image.size
    image.thumbnail((math.ceil(w/4), math.ceil(h/4)))
    return image

def preprocess_image(file_path):
    image = read_image(file_path)
    downsample_image = downsample_image(image)
    return downsample_image

def get_image_data(image):
    return np.array(image)

def list_files(directory):
    return [f for f in listdir(directory) if isfile(join(directory, f))]
```

```
In [45]: subject_1 = []
subject_2 = []
for file in list_files("data/yalefaces/"):
    if file == 'subject01-test.gif' or file == 'subject02-test.gif':
        continue
    im = preprocess_image("data/yalefaces/" + file)
    im_dat = get_image_data(im).flatten()
    if 'subject01' in file:
        subject_1.append(im_dat)
    else:
        subject_2.append(im_dat)
subject_1 = np.vstack(subject_1)
subject_2 = np.vstack(subject_2)
```

```
In [46]: subject01_00=Image.open('data/yalefaces/subject01/glasses.gif')
```

```
In [47]: Subject01_00.size
```

```
Out[47]: (320, 243)
```

```
In [48]: width, height= Subject01_00.size
```

```
In [49]: w, h=math.ceil(width/4), math.ceil(height/4)
```

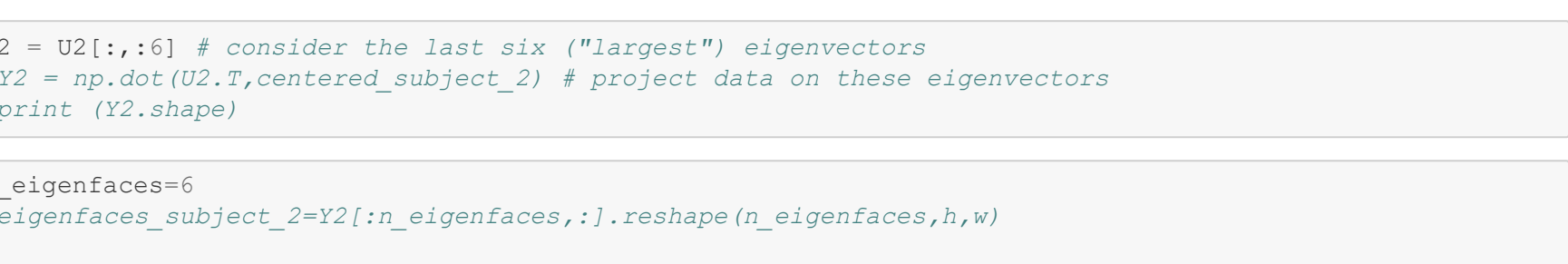
Subject-1 related 10 gifs:

```
In [67]: plt.rcParams["axes.grid"] = False
f, axarr=plt.subplots(1,10,figsize=(20,30))
for i in range(10):
    #plt.subplots(1,1,figsize=(3,4))
    axarr[i].imshow(subject_1[i].reshape((h, w)), cmap = 'gray')
    axarr[i].set_title("Original, number=%i" % (i+1), fontsize=8)
```



Subject-2 related 10 gifs:

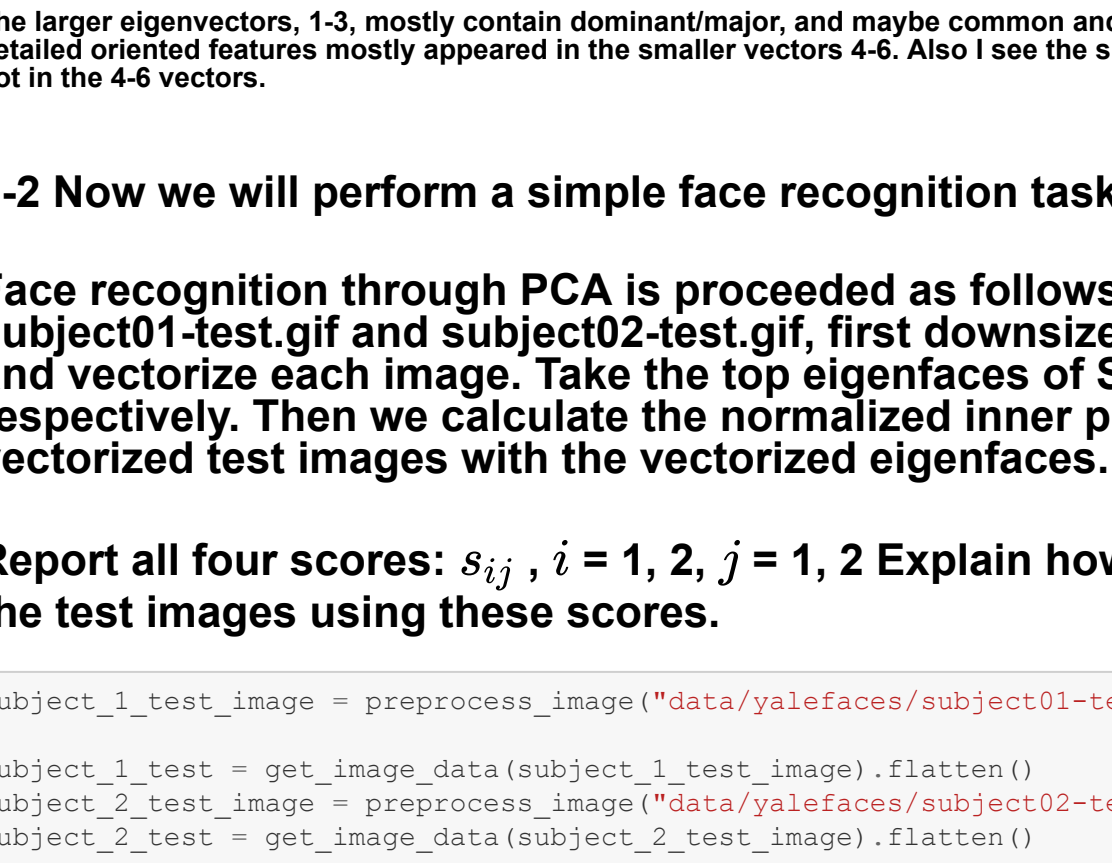
```
In [68]: f, axarr=plt.subplots(1,9,figsize=(20,30))
for i in range(9):
    axarr[i].imshow(subject_2[i].reshape((h, w)), cmap = 'gray')
    axarr[i].set_title("Original, number=%i" % (i+1), fontsize=8)
```



```
In [69]: mean_subject_1 = np.mean(subject_1, axis=0)
mean_subject_2 = np.mean(subject_2, axis=0)
centered_subject_1 = subject_1-mean_subject_1
centered_subject_2 = subject_2-mean_subject_2
```

```
In [70]: f, axarr=plt.subplots(1,2,figsize=(10,15))
axarr[0].imshow(mean_subject_1.reshape((h, w)), cmap = 'gray')
axarr[0].set_title("Mean Subject-1", fontsize=8)
axarr[1].imshow(mean_subject_2.reshape((h, w)), cmap = 'gray')
axarr[1].set_title("Mean Subject-2", fontsize=8)
```

```
Out[70]: Text(0.5, 1.0, 'Mean Subject-2')
```



Calculating the Eigenfaces using Eigen Decomp. :

```
In [71]: C1 = np.cov(centered_subject_1, rowvar=False)
```

```
In [72]: #E1, U1 = np.linalg.eigh(C1)
E1, U1 = np.linalg.eig(C1)
B1, U1 = E1.real / U1.real
```

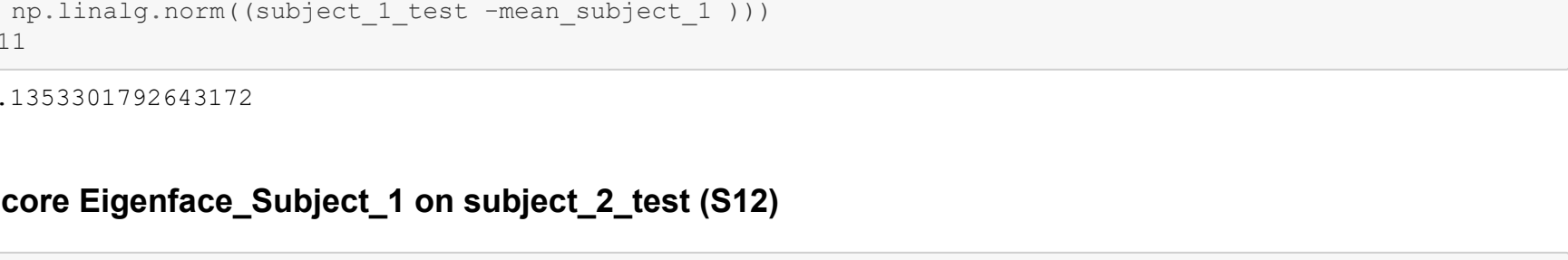
```
In [73]: idc1 = np.argsort(-B1)
#print(idc1[:9]) # first 9 indices Descending Order
#print(B1[idc1[:9]]) # first 9 sorted eigenvalues
```

```
In [74]: B1 = B1[idc1]
U1 = U1[:,idc1] # sort columns
```

```
In [75]: U1 = U1[:,16] # consider the last six largest eigenvectors
np.dot(U2.T, centered_subject_1) # project data on these eigenvectors
print (U1.shape)
```

```
(4880, 6)
```

```
In [76]: n_eigenfaces=6
#eigenfaces_subject_1=Y1[in_eigenfaces,:].reshape(n_eigenfaces,h,w)
f, axarr=plt.subplots(1,n_eigenfaces,figsize=(20,30))
for i in range(n_eigenfaces):
    #axarr[i].imshow(eigenfaces_subject_1[i,:].reshape((h, w)), cmap = 'gray')
    axarr[i].imshow(U1[:,i].reshape((h, w)), cmap = 'gray')
    axarr[i].set_title("Eigenface number=%i" % (i+1))
```



```
In [77]: ### sanity check of eigenvectors U1:
U1[:,0].shape
```

```
Out[77]: (4880,)
```

```
In [78]: U1[:,0].T@U1[:,1]
```

```
Out[78]: 6.192962809237201e-16
```

```
In [79]: U1[:,0].T@U1[:,0]
```

```
Out[79]: 1.0
```

```
In [80]: ### Sanity check was successful for U1
```

```
In [81]: C2 = np.cov(centered_subject_2, rowvar=False)
```

```
In [82]: #E2, U2 = np.linalg.eigh(C2)
E2, U2 = np.linalg.eig(C2)
B2, U2 = E2.real / U2.real
```

```
In [83]: idc2 = np.argsort(-B2)
#print(idc2[:9])
#print(B2[idc2[:9]])
```

```
In [84]: B2 = B2[idc2]
U2 = U2[:,idc2]
```

```
In [85]: U2 = U2[:,16] # consider the last six ("largest") eigenvectors
#U2 = np.dot(U2.T, centered_subject_2) # project data on these eigenvectors
print (U2.shape)
```

```
In [86]: n_eigenfaces=6
#eigenfaces_subject_2=Y2[in_eigenfaces,:].reshape(n_eigenfaces,h,w)
f, axarr=plt.subplots(1,n_eigenfaces,figsize=(20,30))
for i in range(n_eigenfaces):
    #axarr[i].imshow(eigenfaces_subject_2[i,:].reshape((h, w)), cmap = 'gray')
    axarr[i].imshow(U2[:,i].reshape((h, w)), cmap = 'gray')
    axarr[i].set_title("Eigenface number=%i" % (i+1))
```



```
In [87]: ### sanity check eigenvectors U2:
U2[:,0].shape
```

```
Out[87]: (4880,)
```

```
In [88]: U2[:,0].T@U2[:,1]
```

```
Out[88]: -1.9376861226660935e-15
```

```
In [89]: U2[:,0].T@U2[:,0]
```

```
Out[89]: 1.0000000000000002
```

```
In [90]: ### Sanity check was successful for U2
```

Observations on the Eigenfaces:

The larger eigenfactors, 1-3, mostly contain dominant/major, and maybe common and regional, features where as other minor and detailed oriented features mostly appeared in the smaller vectors 4-6. Also I see the shadows only in the first three vectors, and not in the 4-6 vectors.

2-2 Now we will perform a simple face recognition task.

Face recognition through PCA is proceeded as follows. Given the test image subject01-test.gif and subject02-test.gif, first downsize by a factor of 4 (as before), and vectorize each image. Take the top eigenfaces of Subject 1 and Subject 2, respectively. Then we calculate the normalized inner product score of the 2 vectorized test images with the vectorized eigenfaces.

Report all four scores: s_{ij} , $i = 1, 2$, $j = 1, 2$ Explain how to recognize the faces of the test images using these scores.

```
In [91]: subject_1_test_image = preprocess_image("data/yalefaces/subject01-test.gif")
subject_1_test = get_image_data(subject_1_test_image).flatten()
subject_2_test_image = preprocess_image("data/yalefaces/subject02-test.gif")
subject_2_test = get_image_data(subject_2_test_image).flatten()
```

subject_1_test_image size

```
In [92]: subject_1_test.shape
```

```
Out[92]: (4880,)
```

```
In [93]: f, axarr=plt.subplots(1,2,figsize=(10,15))
axarr[0].imshow(subject_1_test.reshape((h, w)), cmap = 'gray')
axarr[0].set_title("Subject-1 Test minus the mean", fontsize=8)
axarr[1].imshow(subject_2_test-mean_subject_2).reshape((h, w)), cmap = 'gray')
axarr[1].set_title("Subject-2 Test minus the mean", fontsize=8)
```

```
Out[93]: Text(0.5, 1.0, 'Subject-2 Test minus the mean')
```



```
In [94]: f, axarr=plt.subplots(1,2,figsize=(10,15))
axarr[0].imshow(subject_1_test-mean_subject_1).reshape((h, w)), cmap = 'gray')
axarr[0].set_title("Subject-1 Test minus the mean", fontsize=8)
axarr[1].imshow(subject_2_test-mean_subject_2).reshape((h, w)), cmap = 'gray')
axarr[1].set_title("Subject-2 Test minus the mean", fontsize=8)
```

```
Out[94]: Text(0.5, 1.0, 'Subject-2 Test minus the mean')
```



```
In [95]: Eigenface_Subject_1 = U1[:,0]
Eigenface_Subject_1.shape
```

```
Out[95]: (4880,)
```

```
In [96]: Eigenface_Subject_2 = U2[:,0]
Eigenface_Subject_2.shape
```

```
Out[96]: (4880,)
```

$S_{ij}^{(k)}$ calculation using the 1st Eigenfaces & Test Images with subtracted Means:

Score Eigenface_Subject_1 on subject_1_test (S11)

```
In [97]: S11=Eigenface_Subject_1.T @ (subject_1_test -mean_subject_1) / (np.linalg.norm(Eigenface_Subject_1.T) * np.linalg.norm(subject_1_test -mean_subject_1))
S11
```

```
Out[97]: 0.1353301792643172
```

Score Eigenface_Subject_1 on subject_2_test (S12)

```
In [98]: S12=Eigenface_Subject_1.T @ (subject_2_test -mean_subject_1) / (np.linalg.norm(Eigenface_Subject_1.T) * np.linalg.norm(subject_2_test -mean_subject_1))
S12
```

```
Out[98]: -0.269154019514262
```

Score Eigenface_Subject_2 on subject_1_test (S21):

```
In [99]: S21=Eigenface_Subject_2.T @ (subject_1_test -mean_subject_2) / (np.linalg.norm(Eigenface_Subject_2.T) * np.linalg.norm(subject_1_test -mean_subject_2))
S21
```

```
Out[99]: -0.560042149671182
```

Score Eigenface_Subject_2 on subject_2_test (S22):

```
In [100]: S22=Eigenface_Subject_2.T @ (subject_2_test -mean_subject_2) / (np.linalg.norm(Eigenface_Subject_2.T) * np.linalg.norm(subject_2_test -mean_subject_2))
S22
```

```
Out[100]: 0.43794152253843166
```

```
In [101]: ##### Sij calculation using the 1st Eigenfaces & Test Images minus the means:
print("S11= (%.2f) , S22=(%.2f)" % (S11,S12))
print("S21= (%.2f) , S22=(%.2f)" % (S21,S22))

S11= 0.14 , S22=0.27
S21= -0.56 , S22=0.44
```

Observations:

I calculated the cosine similarity measure for the test subjects, with mean subtracted, and the largest eigenfaces. I had to multiply -1 into my eigenfaces for calculation of the cosine similarity, just to have all the S numbers make some sense. That may be an acceptable practice as the eigenvectors are not unique, if A is an eigenvector so is the -A, and this may cause a confusion. I had then S11 and S22 as positive numbers and S12 and S21 as negative numbers.

There are other similarity measures, like residual distance calculation that doesn't suffer from this non uniqueness confusion, I implemented that measure in the next part.

2-3 Explain if face recognition can work well and discuss how we can improve it possibly.

We observed in the earlier part the non-uniqueness of the eigenvectors. The class was advised to test another similarity measure, Residual Distance, which does not suffer from this non-uniqueness issue. It simply measure the distance between any two vectors, in our case U is the eigenvector and I is the test images:

Another similarity measure, the Residual Distance Definition:

$$R = (I - (UU^T))^2$$

We calculated the pair residual distances, R :

Those R numbers are always positive, since $(-U)(-U)^T = UU^T$, and hence I believe the residual distance measure may be more effective than the cosine similarity measure.

Residual distance calculation for the EigenFaces and the test images subtracted means

```
In [102]: ###Residual Distance for Eigenvecotor-1 & Test Image-1
t= (subject_1_test-mean_subject_1).reshape((4880,1))
u = Eigenface_Subject_1.reshape((4880,1))
rej1_1=np.linalg.norm(t-u@u.T@t)
```

```
"(t-u).format(rej1_1)
```

```
Out[102]: '2.593065e+03'
```

```
In [103]: ###Residual Distance for Eigenvecotor-1 & Test Image-2
t= (subject_2_test-mean_subject_1).reshape((4880,1))
u = U1[:,0].reshape((4880,1))
rej1_2=np.linalg.norm(t-u@u.T@t)
```

```
"(t-u).format(rej1_2)
```

```
Out[103]: '6.315078e+03'
```

```
In [104]: ###Residual Distance for Eigenvecotor-2 & Test Image-1
t= (subject_1_test-mean_subject_2).reshape((4880,1))
u = Eigenface_Subject_2.reshape((4880,1))
rej2_1=np.linalg.norm(t-u@u.T@t)
```

```
"(t-u).format(np.linalg.norm(rej2_1))
```

```
Out[104]: '5.745946e+03'
```

```
In [105]: ###Residual Distance for Eigenvecotor-2 & Test Image-2
t= (subject_2_test-mean_subject_2).reshape((4880,1))
u = Eigenface_Subject_2.reshape((4880,1))
rej2_2=np.linalg.norm(t-u@u.T@t)
```

```
"(t-u).format(np.linalg.norm(rej2_2))
```

```
Out[105]: '1.915918e+03'
```

```
In [106]: #####Residual Distance calculation using the 1st Eigenfaces & Test Images-means:
print("Ultimate Residual Distance 1_1= (%.2f) , Ultimate Residual Distance 1_2= (%.2f)" % (min(M_D11),min(M_D12)))
print("Ultimate Residual Distance 2_1= (%.2f) , Ultimate Residual Distance 2_2= (%.2f)" % (min(M_D21),min(M_D22)))

Ultimate Residual Distance 1_1= 2593.07 , Ultimate Residual Distance 1_2= 6315.08
Ultimate Residual Distance 2_1= 5745.95 , Ultimate Residual Distance 2_2= 1915.92
```

The above number now make sense, and there is also no confusion regarding the non-uniqueness.

I also extended this measure to the other EigenVectors, I had six of them for each image, and investigated if any of the eigenvector could further minimize such Distance:

```
In [107]: def Min_residual_distance_all(eigenfaces, subject,k):
    subject_r=subject.reshape((4880,1))
    distance =[]
    sum_distance=0
    for i in range(0,k):
        eigenface_temp =eigenfaces[:,i].reshape((4880,1))
        #eigenface_temp =eigenfaces[:,i]
```

```
#sum_distance = (np.linalg.norm((I-(eigenface_temp@Eigenface_temp.T))@subject_r))
sum_distance = (np.linalg.norm(subject_r -eigenface_temp@Eigenface_temp.T @ subject_r))
distance[i]= np.linalg.norm(sum_distance)
```

```
return sorted(distance.values())
```

```
In [108]: M_D11 = Min_residual_distance_all(U1, (subject_1_test-mean_subject_1),6)
M_D11
```

```
Out[108]: [2507.7625956707757,
2550.3144922676306,
2593.06315083437,
2601.1557767442546,
2603.970173752224,
2611.3291033537066]
```

```
In [109]: M_D12 = Min_residual_distance_all(U1, (subject_2_test-mean_subject_1),6)
M_D12
```

```
Out[109]: [6206.84725069171,
6315.0784362187305,
6518.615718319646,
6531.045832609504,
6531.8804993627685,
6540.185714466218]
```

```
In [110]: M_D21 = Min_residual_distance_all(U2, (subject_1_test-mean_subject_2),6)
M_D21
```

```
Out[110]: [5745.946161748096,
6830.395143230939,
6901.76103044426,
6918.60528771121,
6919.353618522969,
6933.847073875485]
```

```
In [111]: M_D22 = Min_residual_distance_all(U2, (subject_2_test-mean_subject_2),6)
M_D22
```

```
Out[111]: [1751.9347243767537,
1715.9176833084649,
2074.2981949208353,
2120.5050889453278,
2121.501775819062,
2129.5074063258607]
```

```
In [112]: #####Residual Distance calculation using All the Eigenfaces & Test Images-means:
print("Ultimate Residual Distance 1_1= (%.2f) , Ultimate Residual Distance 1_2= (%.2f)" % (min(M_D11),min(M_D12)))
print("Ultimate Residual Distance 2_1= (%.2f) , Ultimate Residual Distance 2_2= (%.2f)" % (min(M_D21),min(M_D22)))

Ultimate Residual Distance 1_1= 2507.76 , Ultimate Residual Distance 1_2= 6206.85
Ultimate Residual Distance 2_1= 5745.95 , Ultimate Residual Distance 2_2= 1915.92
```

Conclusion:

The above numbers provide further confidence in such face recognition practice. I looked in all the eigenvectors and observed that the minimum distance numbers are associated with either 1st or 2nd eigenface. The results also make sense in linking the relevant eigenface & test-image.