

## HW7 REPORT

### ● PCA project

#### The PCA class

```
4 class PCA:
5     def __init__(self, dataArray, classArray=None):
6         #flattens array to 2-D
7         if dataArray.ndim>2:
8             dataArray=dataArray.reshape(dataArray.shape[0],-1)
9         self.dataArray=dataArray
10        self.classArray = classArray
11        self.dataSize, self.dataDimension = self.dataArray.shape
```

`__init__`:

將資料壓成 [資料比數 , feature 數]，正確的 class 可有可無，如果有給的話最後的圖會上色。定義一些基本的變數。

```
12 def project(self, projectDimension=2):
13     #project process here
14     self.projectDimension=projectDimension
15
16     print("Calculating covariance matrix...")
17     print(self.dataArray.shape)
18     covarArray=np.cov(self.dataArray, rowvar=False)
19     #covarArray=np.corrcoef(self.dataArray, rowvar=False)
20
21     print("Solving Eigen Values...")
22     self.eigenValue, self.eigenVector = linalg.eig(covarArray)
23     #reverse argsort() indexes (was sorted in ascending)
24     sortIndices=self.eigenValue.argsort()[::-1][:]
25     self.eigenValue=self.eigenValue[sortIndices]
26     self.eigenVector = self.eigenVector[:, sortIndices]
27
28     print("Projecting data...")
29     #select eigenvectors with largest eigen values
30     self.projectEigenVector=self.eigenVector[:, :self.projectDimension]
31     self.projectData=self.dataArray@self.projectEigenVector
32     print("Finished projecting!")
```

`project()`:

將資料投射到 `projectDimension` 的維度空間。covariance 因為自己用迴圈跑很慢因此借用 `np.cov`。算 eigenvalues 用的是 numpy 的 `linalg`。 `argsort()` 回傳的是 ascending 的 index order，我要取最大值所以將其反向。Project 就是矩陣乘法乘剛剛的 eigenvectors，這邊因為方向的關係可能跟常見公式不同但算法是一樣的。

```

34 def plot(self):
35     #plt figure
36     fig = plt.figure()
37     ax = fig.add_subplot(111)
38     #draw corresponding first two eigenvectors values..
39     dataX=self.projectData[:,0]
40     dataY=self.projectData[:,1]
41     scatter = ax.scatter(dataX, dataY, c=self.classArray, s=20)
42     ax.set_xlabel('x')
43     ax.set_ylabel('y')
44     #classArray can be empty (it's unsupervised after all)
45     if self.classArray is not None:
46         plt.colorbar(scatter)
47     plt.show()

```

plot():

使用 matplotlib 做圖，只取前兩個值作圖，因此若 `projectDimension>2` 會沒辦法呈現後面的值。如果有給各個資料的 `class(label)` 的話會上色。

## PCA PROJECT

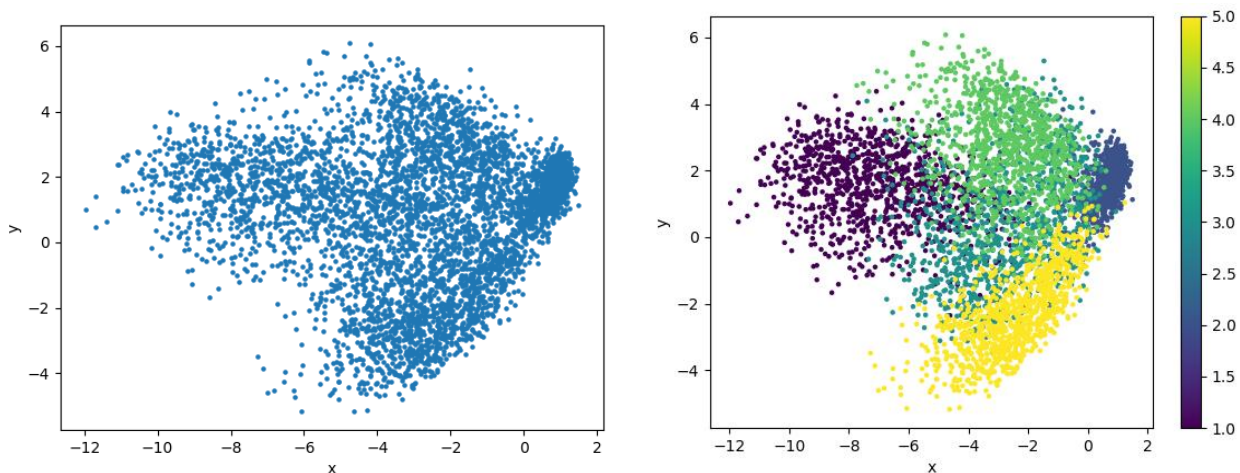
```

48 if __name__=="__main__":
49     from dataprocess import readData
50     trainData=readData("X_train.csv")
51     dataLabel=readData("T_train.csv")
52     pcaTest=PCA(trainData)
53     pcaTest.project()
54     pcaTest.plot()

5 def readData(name,delimiter=','):
6     #load csv file to numpy array
7     print("Reading data {}".format(name))
8     data=np.loadtxt(name,delimiter=delimiter)
9     if data.ndim==2:
10         return data[:,:]
11     elif data.ndim==1:
12         return data[:]

```

讀檔後丟進去。



右邊的上色是根據正確的 `class` 去做的，可以觀察到雖然有大面積的重疊，但是可以大致上看出各個數字有各自的區塊，將 784 壓至 2 有這樣的結果已經算是不錯

了。如果用一點想像去看，可以觀察到 235 這幾個比較像的數字的交疊在一起，而 4 這個與其他的完全不像的數字則有很小的 variance。

## ● LDA Project

### The LDA class

```
4 class LDA:
5     def __init__(self, dataArray, classArray):
6         self.dataArray = dataArray
7         self.classArray = classArray
8         self.realClasses, self.classDataCount = np.unique(self.classArray, return_counts=True)
9         self.totalClassCount = self.realClasses.size
10        self.classDict = dict(zip(self.realClasses, range(self.totalClassCount)))
11        self.dataSize, self.dataDimension = self.dataArray.shape

__init__():
```

LDA 是 supervised learning，所以 data 跟 class 都是必要的。np.unique 是為了找出 class 的種類，classDict 則給在 array 中他對照的 index。

```
13 def project(self, projectDimension=2):
14     self.projectDimension = projectDimension
15     self.inClassScatter = np.zeros((self.dataDimension, self.dataDimension))
16     self.betweenClassScatter = np.zeros((self.dataDimension, self.dataDimension))
17     #calculate class means
18     print("Calculating in class scatter...")
19     inClassMean = np.zeros((self.dataDimension, self.totalClassCount))
20     for dataIndex in range(self.dataSize):
21         inClassMean[:, self.classDict[self.classArray[dataIndex]]] += self.dataArray[dataIndex, :]
22     for classIndex in range(self.totalClassCount):
23         inClassMean[:, classIndex] /= self.classDataCount[classIndex]
24     #calculate inclass scatter (sum over in class)
25     for dataIndex in range(self.dataSize):
26         diffToClassMean = (self.dataArray[dataIndex, :] -
27                             inClassMean[:, self.classDict[self.classArray[dataIndex]]])[:, np.newaxis]
28         self.inClassScatter += (diffToClassMean @ np.transpose(diffToClassMean))
29
30     #calculate between class scatter
31     print("Calculate between class scatter..")
32     allMean = np.average(inClassMean, axis=1, weights=self.classDataCount)[:, np.newaxis]
33     for classIndex in range(self.totalClassCount):
34         diffToTotalMean = inClassMean[:, classIndex] - allMean
35         self.betweenClassScatter = self.classDataCount[classIndex] * (diffToTotalMean @ np.transpose(diffToTotalMean))
36
37     #calculate eigen vectors
38     print("Calculate eigens...")
39     inverseSwSb = linalg.pinv(self.inClassScatter) @ self.betweenClassScatter
40     self.eigenValue, self.eigenVector = linalg.eig(inverseSwSb)
41     sortIndices = self.eigenValue.argsort()[::-1][:self.projectDimension]
42     self.eigenValue = self.eigenValue[sortIndices]
43     self.eigenVector = self.eigenVector[:, sortIndices]
44     print(self.eigenVector)
45     self.projectData = self.dataArray @ self.eigenVector
```

project():

這整塊簡單說就是 算 inclass scatter -> between class scatter -> eigenvectors -> project

inclass mean 是指投影片中的 mj，接著將各個點和對應的 class mean 相減後加總

就是 in class scatter。Self.classDict[] 就是為了對應實際 class 跟存在 array 中的位置

而已。有些行數有加上[:,np.newaxis]，這是幫 1D 加上一個軸，這樣才能做矩陣乘法@。

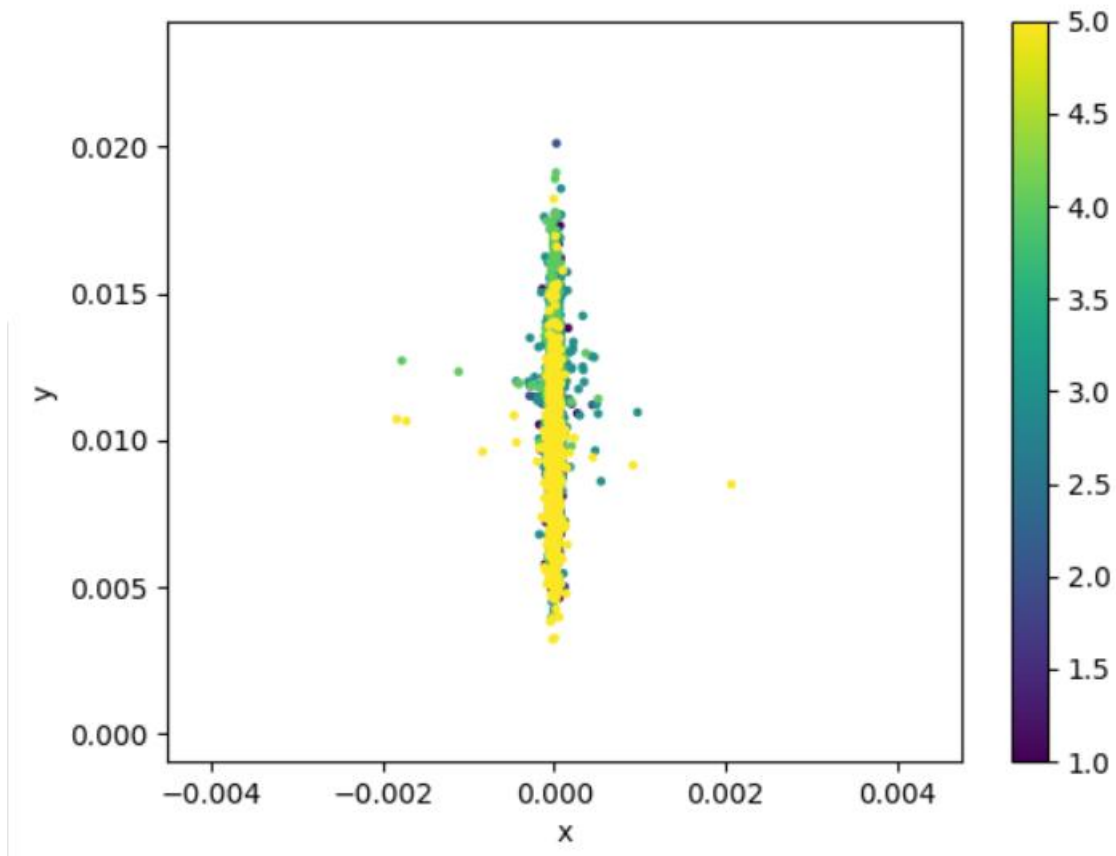
再來就是算出投射用的 eigenvector，是指 $(\text{inclass scatter})^{-1} * (\text{between class scatter})$ ，在此就不推導。因為 inclass scatter 有可能沒有 inverse，因此使用 pinv(pseudo inverse)。最後取 eigenvalue 最大的 eigenvector 做 project。

```
46 def plot(self):
47     fig = plt.figure()
48     ax = fig.add_subplot(111)
49     dataX=self.projectData[:,0]
50     dataY=self.projectData[:,1]
51     scatter = ax.scatter(dataX, dataY, c=self.classArray, s=5)
52     ax.set_xlabel('x')
53     ax.set_ylabel('y')
54     plt.colorbar(scatter)
55     plt.show()
```

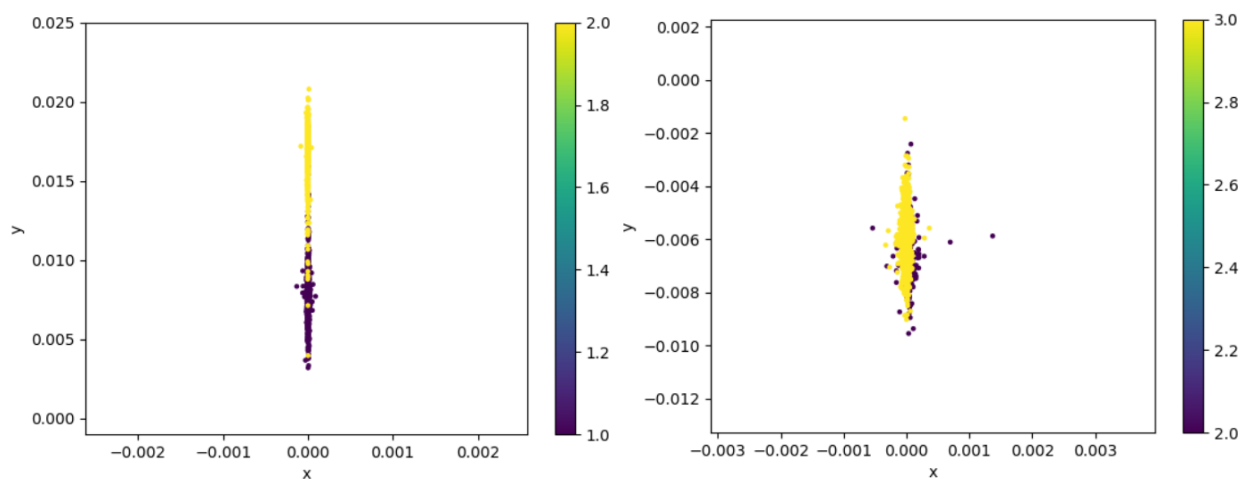
和 PCA 一樣的 plot。

## LDA PROJECT

```
57 if __name__=="__main__":
58     from dataprocess import readData
59     trainData=readData("X_train.csv")
60     dataLabel=readData("T_train.csv")
61     ldaTest=LDA(trainData,dataLabel)
62     ldaTest.project()
63     ldaTest.plot()
```



這邊的 LDA 分布會都重疊，乍看之下很失敗，但若一次只用兩個 class 可以發現事實上是有。



左圖是只用 1,2 class，可以觀察到有明顯的分開，而右圖 2,3 class 在 PCA 中可以看到本來就相似，因此幾乎無法分開，這算是預期的結果。

## ● Eigen face

```

6 def plotSampleFace(imageArray, imageShape=None, sampleShape=(5,5), randomSelect=True, imageIndex=None, scaling=False):
7     sampleHeight, sampleWidth = sampleShape
8     if imageArray.ndim==2:
9         imageHeight, imageWidth = imageShape
10        imageCount = imageArray[0]
11    elif imageArray.ndim==3:
12        imageCount, imageHeight, imageWidth = imageArray.shape
13    else:
14        print("Illegal image array")
15        return
16    mergedFaces = np.zeros((sampleHeight*imageHeight, sampleWidth*imageWidth))
17    if randomSelect:
18        selectedImageIndex = np.random.randint(imageCount, size=sampleShape)
19    else:
20        if imageIndex is None:
21            selectedImageIndex = np.arange(sampleHeight*sampleWidth).reshape(sampleShape)
22        else:
23            selectedImageIndex = imageIndex
24
25    if scaling is True:
26        for imageIndex in selectedImageIndex:
27            if imageArray.ndim==2:
28                imageArray[imageIndex, :] *= 255.0/imageArray[imageIndex, :].max()
29            elif imageArray.ndim==3:
30                imageArray[imageIndex, :, :] *= 255.0 / imageArray[imageIndex, :, :].max()
31    for sampleRow in range(sampleHeight):
32        if imageArray.ndim==2:
33            mergedImageRow = mergeImageToRow(imageArray[selectedImageIndex[sampleRow, :], :], (imageHeight, imageWidth))
34        elif imageArray.ndim==3:
35            mergedImageRow = mergeImageToRow(imageArray[selectedImageIndex[sampleRow, :, :], :], (imageHeight, imageWidth))
36        mergedFaces[sampleRow*imageHeight:(sampleRow+1)*imageHeight, :] = mergedImageRow
37
38    plt.imshow(mergedFaces)
39    plt.colorbar()

```

## plotSampleFace():

randomSelect 決定 index 是給定，按順序，或是隨機。Scaling 是將值移到 0~255 間，此處沒有用到。

這個函式最主要是要將影像的 row 接起來，其中 mergeImageToRow() 就是負責將一行影像接起來的函式，我再將好幾行疊起來。

```

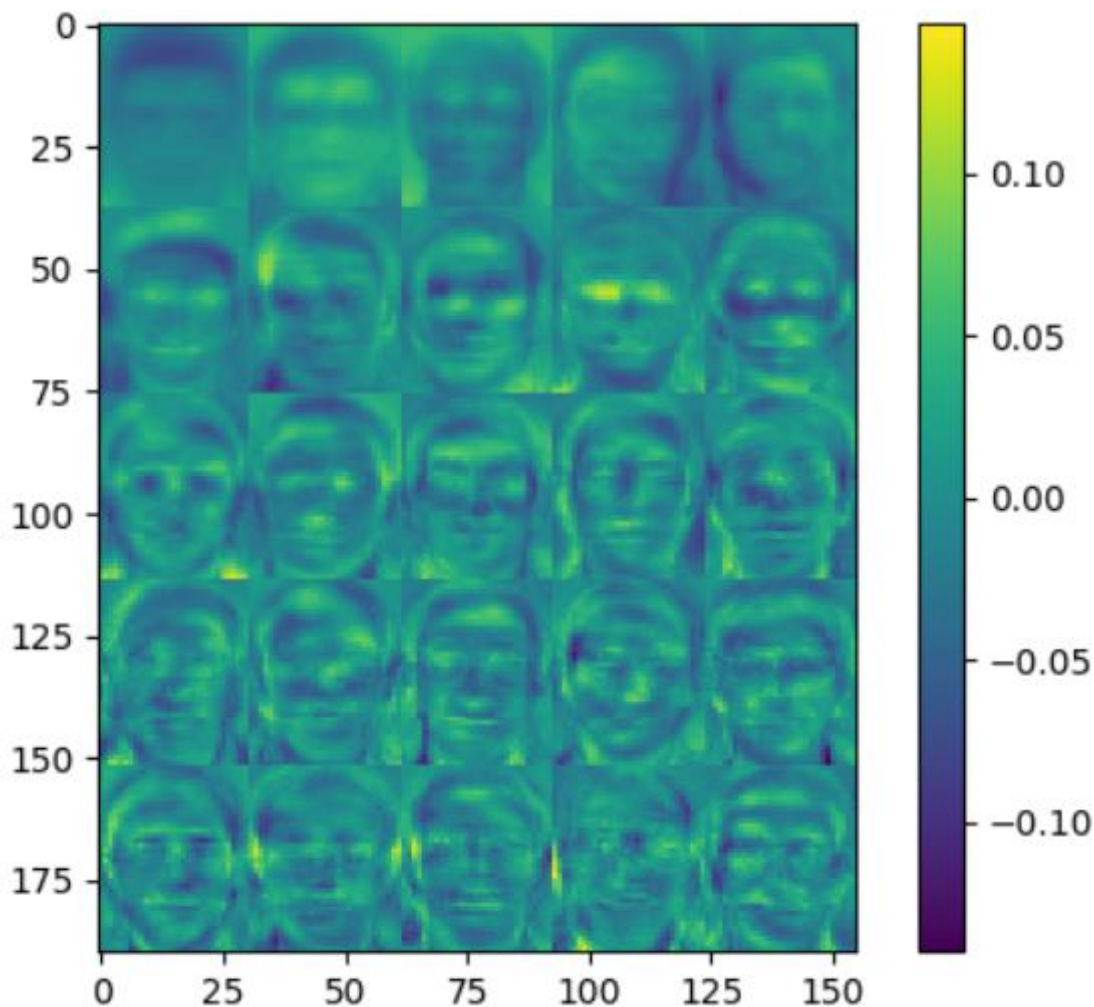
43 def mergeImageToRow(imageArray, imageShape=None):
44     #print(imageShape)
45     #plt.imshow(imageArray[0, :, :])
46     #plt.show()
47     if imageArray.ndim==2:
48         imageHeight, imageWidth = imageShape
49         imageCount = imageArray.shape[0]
50     elif imageArray.ndim==3:
51         imageCount, imageHeight, imageWidth = imageArray.shape
52     else:
53         print("Illegal image array")
54         return
55     mergedImage = np.zeros((imageHeight, imageWidth*imageCount))
56     for imageRow in range(imageHeight):
57         for imageIndex in range(imageCount):
58             if imageArray.ndim==2:
59                 mergedImage[imageRow, imageIndex*imageWidth:(imageIndex + 1) * imageWidth] = imageArray[imageIndex, imageRow, :]
60             elif imageArray.ndim==3:
61                 mergedImage[imageRow, imageIndex*imageWidth:(imageIndex + 1) * imageWidth] = imageArray[imageIndex, imageRow, :]
62     return mergedImage

```

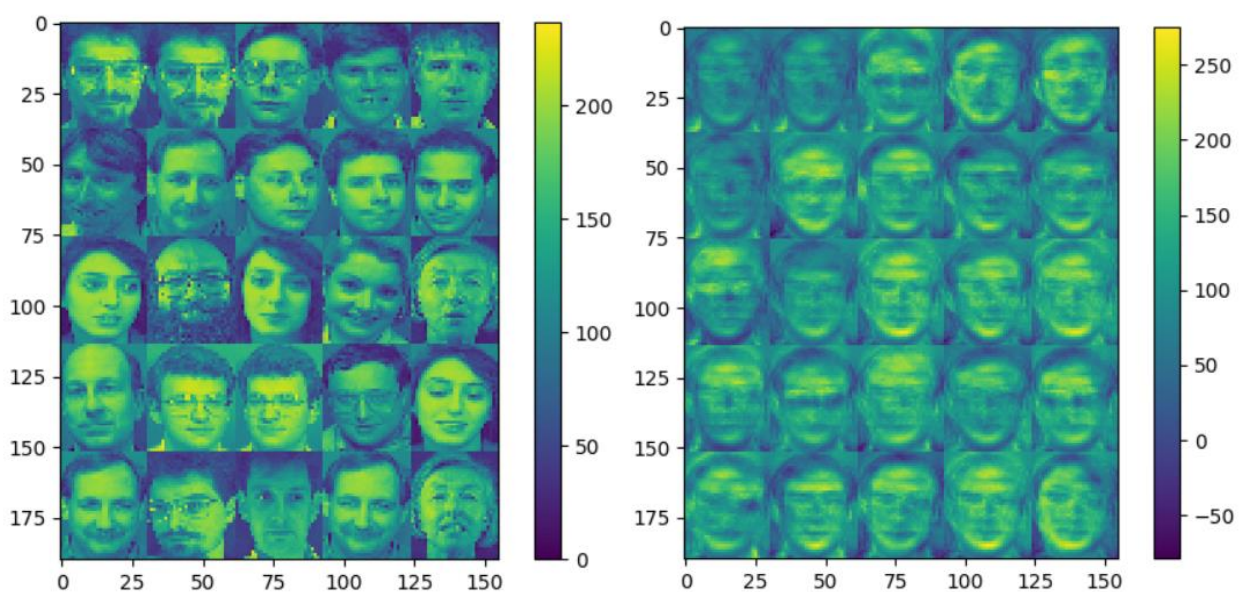
## mergeImageToRow():

將一行影像接起來，以這次實驗為例就是輸出 1x5 的影像數。

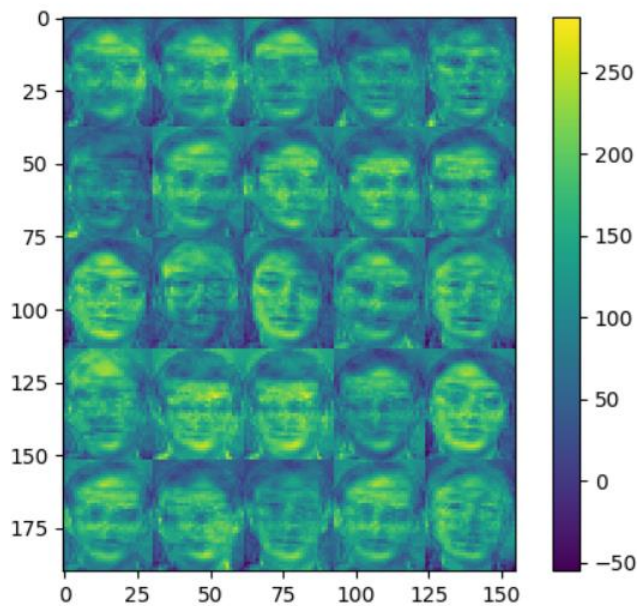




這是根據 400 個臉做出的 eigenface，從左上到右下是 eigenface 的順序。觀察得到左上比較模糊，這是比重比較高比較通用的 eigenface，細節也會比較少。越是到右下角的細節就越多，是對於對應特定圖片細節的 face。



左邊是原圖，右邊是 25 eigenvector 的重構。



這是 50 eigenvector 的重構。可以觀察到 25 eigenvector 的重構各圖片比較相似，50 eigenvector 可以明顯地看到原圖的樣子。

## ● RATIO and NORMALIZED CUT

```

28 #Help with scipy to speed up the calculation
29 if kernel=='linear':
30     self.graphMatrix=linear_kernel(self.dataArray)
31 elif kernel=='RBF':
32     self.graphMatrix = rbf_kernel(self.dataArray, gamma=gamma)
33     #self.graphMatrix=np.maximum(self.graphMatrix-(1e-4),0)
34 elif kernel=='linearRBF':
35     self.graphMatrix=linear_kernel(self.dataArray)+rbf_kernel(self.dataArray, gamma=gamma)

```

用了 scipy 以加速 similarity matrix 的運算。

Ratio 需要 generalized eigenvector

```

65 if mode=='ratio':
66     self.eigenValue, self.eigenVector = linalg.eig(self.graphLaplacian, self.degreeMatrix)
67 else:
68     self.eigenValue, self.eigenVector= linalg.eig(self.graphLaplacian)

```

Normalized cut 需要 normalized laplacian 以及對 row 做 normalize



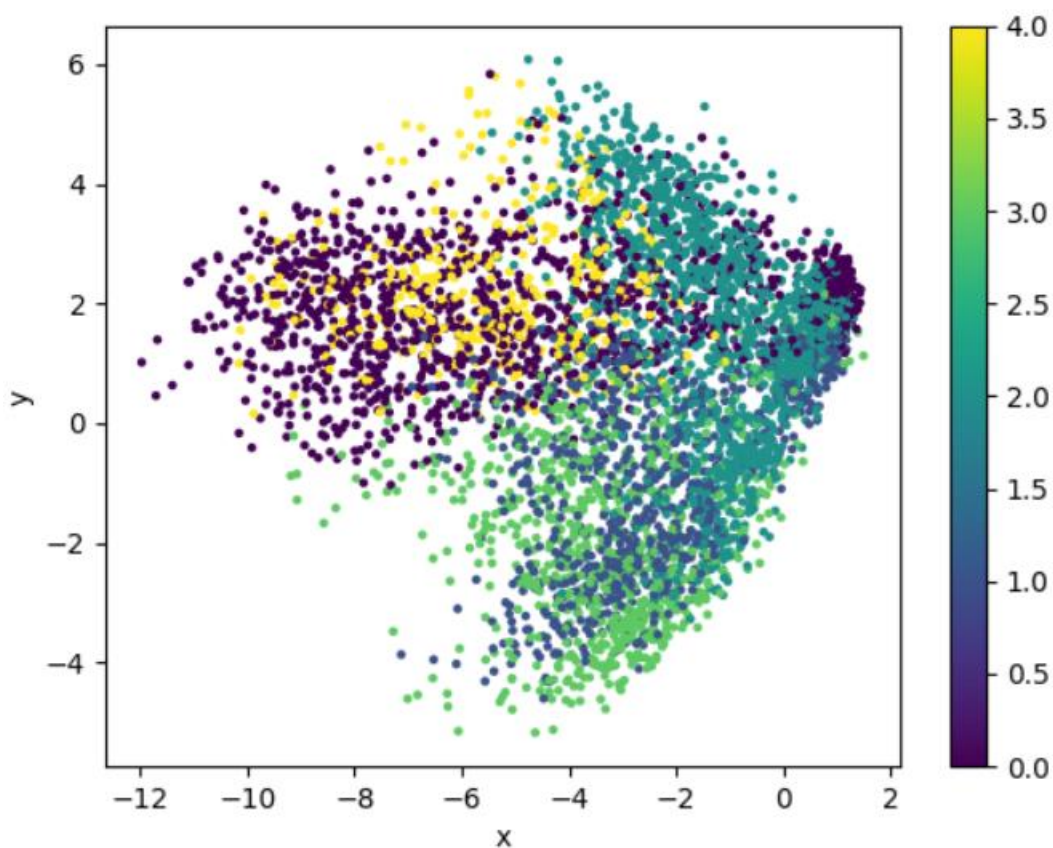
```

55     if mode=='ratio':
56         self.graphLaplacian=self.degreeMatrix-self.graphMatrix
57     elif mode=='normalized':
58         self.graphLaplacian = self.degreeMatrix - self.graphMatrix
59         degreeMatrixInvSqrt=sqrtm(lng.inv(self.degreeMatrix))
60         self.graphLaplacian=degreeMatrixInvSqrt@self.graphLaplacian@degreeMatrixInvSqrt
61     else:
62         print("unknown mode")
63         return

74     if mode=='normalized':
75         for eigenRow in range(self.eigenVector.shape[0]):
76             normTerm=(np.square(self.eigenVector[eigenRow,:]))*0.5
77             self.eigenVector[eigenRow,:]=self.eigenVector[eigenRow,:]/normTerm

```

## ● Normalized\_Linear



## ● Normalized\_Linear+RBF

