

# 2018 ML HW7 Clustering

0756110 李東霖

This document made by HackMD, you can view here <https://hackmd.io/s/Byskl4PeN>  
(<https://hackmd.io/s/Byskl4PeN>)

- 2018 ML HW7 Clustering
  - implementation of Kernels
  - PCA & Kernel PCA
    - kernel trick
    - visualization
    - my observation
  - LDA & Kernel LDA
    - kernel trick
    - visualization
    - my observation
  - RatioCut & NormalizeCut
    - visualization (3 Kernel with 2 cut)
    - my observation
  - SVC result visualization
    - In PCA space
    - In LDA space
    - my observation
  - eigenfaces from att\_faces
    - my observation

## implementation of Kernels

總共實作了四種 kernel :

- Linear
- Polynomial
- RBF
- Linear+RBF

```

1  def __Normal(x):
2      return (x - np.min(x))/(np.max(x) - np.min(x))
3
4  def Linear(u,v):
5      return np.matmul(u, v.T);
6
7  def Poly(u,v, gamma, coef0, degree):
8      return ((gamma*np.matmul(u, v.T)) + coef0)**degree
9
10 def Euclidean(x,y):
11     """
12     calculate Euclidean distance
13     parameters:
14         x: n1 * d, y: n2 x d
15     return:
16         d: n1 * n2
17
18     d(i, j) = |x(i) - y(j)|^2
19     """
20     if len(x.shape)==1:
21         x = x[None,:]
22     if len(y.shape)==1:
23         y = y[None,:]
24     return np.matmul(x**2, np.ones((x.shape[1],y.shape[0]))) \
25         + np.matmul(np.ones((x.shape[0],x.shape[1])), (y**2).T) \
26         - 2*np.dot(x,y.T)
27
28 def RBF(u,v,gamma=1):
29     return np.exp(-1*gamma*Euclidean(u,v));
30
31 def LinearKernel():
32     return lambda u,v:Linear(u,v)
33
34 def PolyKernel(gamma=1, coef0=1, degree=3):
35     return lambda u,v:Poly(u,v,gamma,coef0,degree)
36
37 def RBFKernel(gamma=1):
38     """
39     generate callable function for RBF kernel
40     parameters:
41         gamma : default is 1
42     return:
43         lambda(u, v)
44
45         rbf(u, v) = exp(-gamma|u-v|^2)
46     """
47     return lambda u,v:RBF(u,v,gamma)
48
49 def LinearRBFKernel(gamma=1):
50     return lambda u,v:__Normal(Linear(u,v)) + RBF(u,v,gamma)

```

在 Linear+RBF kernel 中稍微做了點處理

因為 RBF kernel 的數值只會介在 (0, 1] 之間，但是 Linear kernel 卻不是

如果直接相加會導致結果幾乎是等於 Linear kernel

因此先將 Linear kernel 線性變換到 [0, 1] 再進行相加

在之後的使用可以看到 Linear+RBF 結果介在 Linear kernel 與 RBF kernel 之間

這邊設定好幾個 kernel 的參數，方便之後的重複使用與比較

```
51 # parameters setting from HW5 grid search + cross-validation
52 __kernels = {
53     'Linear' : LinearKernel(),
54     'Poly' : PolyKernel(3/1000, 1, 10),
55     'RBF' : RBFKernel(0.01),
56     'Linear+RBF' : LinearRBFKernel(0.01)
57 }
```

## PCA & Kernel PCA

PCA 希望能找到一組投影向量能讓投影過去之後投影回來與原始資料的 MSE 最小

$$\begin{aligned} z &= xW \\ MSE &= ||x - zW^T||^2 = ||x - xWW^T||^2 \\ &= \text{constant} - \text{constant}W^T SW \\ S_{ij} &= \frac{\sum_{k=1}^N (x_k(i) - \bar{x(i)})(x_k(j) - \bar{x(j)})^T}{N}, x = \begin{bmatrix} x(0) \\ x(1) \\ \vdots \\ x(d) \end{bmatrix} \end{aligned}$$

$S$  為 covariance matrix，經由上式得知要找到最小的 MSE 就是找到最大的  $W^T SW$  因此去解出 eigen problem  $SW = \lambda W$ ，並取出擁有最大 eigen value 的  $k$  個 vector 就是 PCA 想找的投影向量也是該份資料的主成份

### kernel trick

除了從原始資料空間進行投影，也可以嘗試從 feature space 去投影

但因為 data space 變換 feature space 不一定容易計算，因此使用 kernel trick 來計算

$$\begin{aligned} \frac{(\phi(X) - \frac{\sum_i^n \phi(X_i)}{N})(\phi(X) - \frac{\sum_i^n \phi(X_i)}{N})^T}{N} W &= \lambda W, W = \phi(X)\alpha \\ \Rightarrow (K - 1_N K - K 1_N + 1_N K 1_N)\alpha &= \lambda N \alpha \\ z &= \phi(x)W = \phi(x)\phi(X)\alpha = K(x, X)\alpha \end{aligned}$$

以下為實作的程式碼，eigen problem 使用 `numpy.linalg.eigh` 解

```

1  def covariance(x, y):
2      x1 = x - (np.sum(x, axis = 1) / x.shape[1])[:,None]
3      y1 = y - (np.sum(y, axis = 1) / y.shape[1])[:,None]
4      return np.matmul(x1, y1.T)/ x1.shape[1]
5
6  def PCA(datas, k, kernel=None):
7      if callable(kernel):
8          datas = kernel(datas, datas)
9          N = datas.shape[0]
10         N1 = np.ones((N,N))/N
11         S = (datas
12             - np.matmul(N1, datas)
13             - np.matmul(datas, N1)
14             + np.matmul(np.matmul(N1, datas), N1)
15             )/N
16     else:
17         S = covariance(datas.T, datas.T)
18
19     value, vector = np.linalg.eigh(S)
20
21     start_idx = 0
22     max_idx = np.flip(np.argsort(value))
23     W = np.concatenate([
24         vector[:, max_idx[i]][:,None]
25         for i in range(start_idx,k+start_idx)
26     ], axis=1)
27
28     # whitening
29     if callable(kernel):
30         W /= np.sqrt(value[max_idx[start_idx:start_idx+k]])[None, :]
31
32     return np.matmul(datas, W), W, value, vector

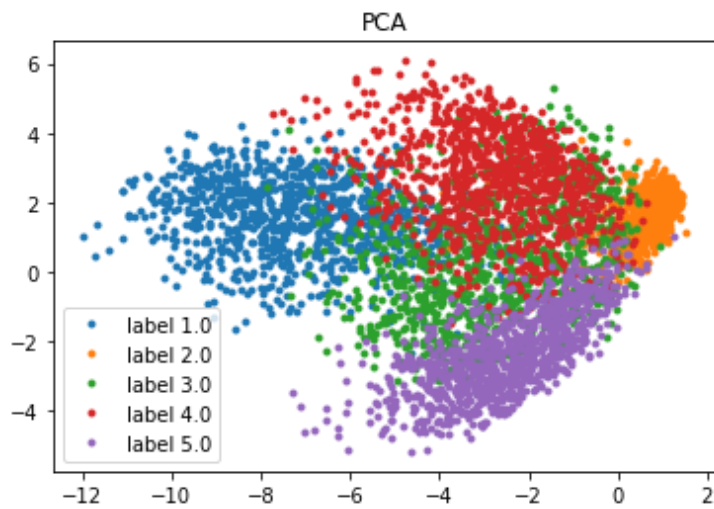
```

回傳的第一個參數即是投影完的空間

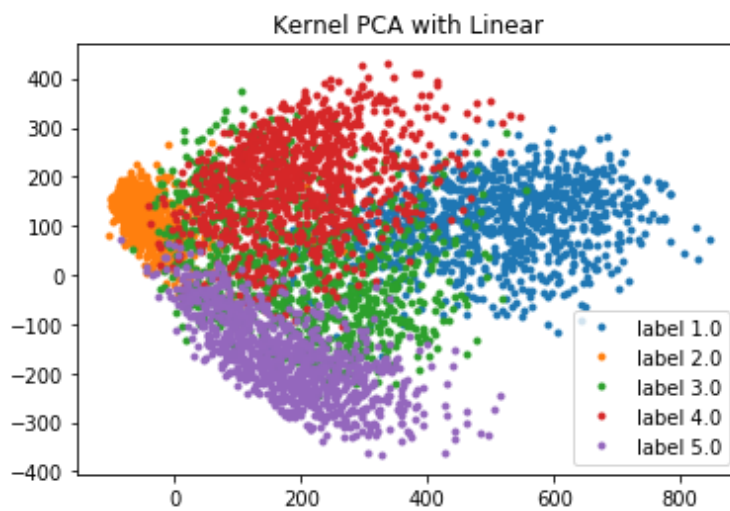
使用 kernel PCA 要注意不是直接將計算得到的 eigen vector 使用在原本資料空間  
而是要使用在 kernel matrix 上，因為實際上取得的 eigen vector 代表的是係數 *alpha*  
這些係數乘上在 feature space 的所有資料點才會是投影向量 *W*

## visualization

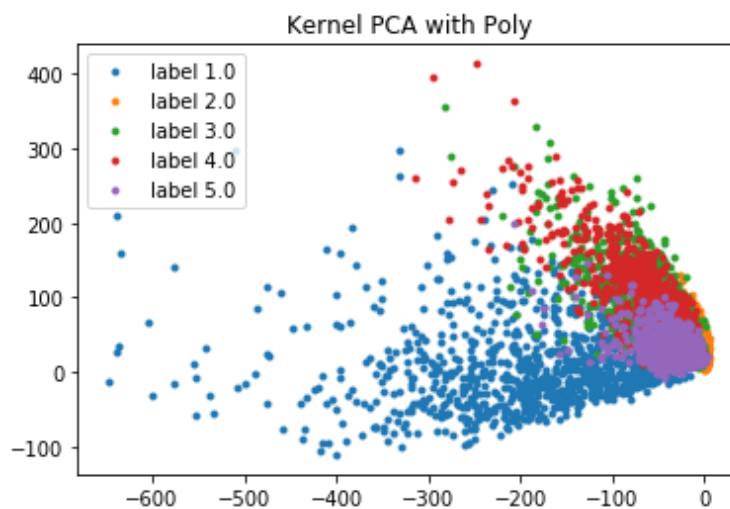
- PCA



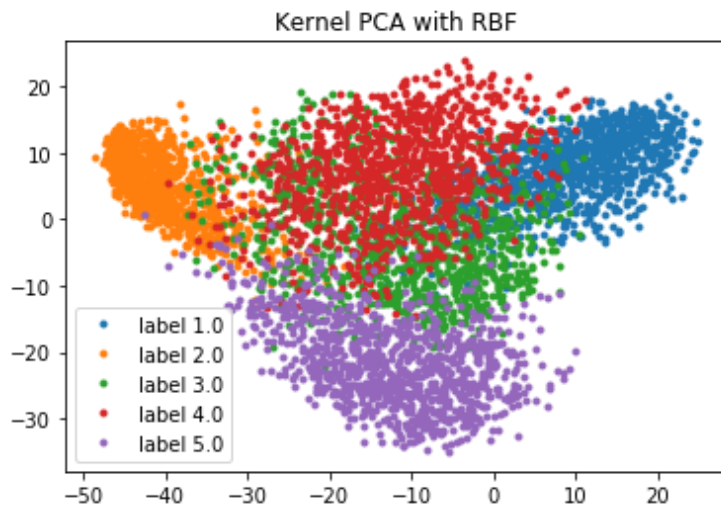
- PCA with Linear kernel



- PCA with Polynomial kernel  
 $\gamma = 3/1000$ ,  $\text{coef0} = 1$ ,  $\text{degree} = 10$

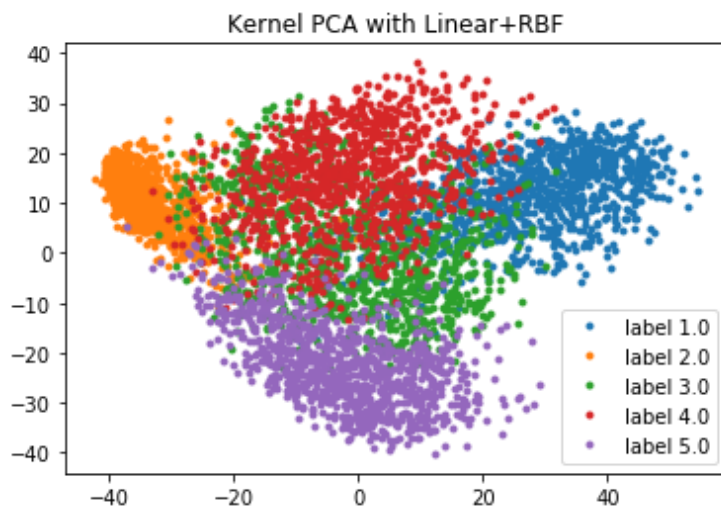


- PCA with RBF kernel  
 $\gamma = 1/100$



- PCA with Linear+RBF kernel

$$\gamma = 1/100$$



### my observation

這邊可以發現 PCA 與 PCA with Linear kernel 是非常相似  
只是形狀上下左右顛倒，這是因為在 Linear kernel 中其實 feature space 與 data space 是一樣的

因此做出來會有類似結果也是合情合理

但畢竟計算過程與維度並不一樣因此座標軸範圍有差距

可以看到 RBF 的結果還不錯，同時 Linear+RBF 則確實介在 Linear 與 RBF 的結果

## LDA & Kernel LDA

LDA 是一種 supervised learning，需要告知資料的標籤

求解過程也容易理解，希望找到一組投影向量  $W$  滿足以下條件

- 使得同一群的中心 (mean) 離所有資料的中心盡可能遠
  - maximize between-class-scatter  $S_B$
- 同一群自己的變異數 variance 盡量小，也就是同群要集中
  - minimize within-class-scatter  $S_W$

如何求解，如下

$$S_B = \sum_{j=1}^k |C_j|(m_j - m)(m_j - m)^T, m = \frac{\sum x}{n}, m_j = \frac{\sum_{x_i \in C_j} x_i}{|C_j|}$$

$$S_W = \sum_{j=1}^k \sum_{x_i \in C_j} (x_i - m_j)(x_i - m_j)^T$$

$$S_B W = \lambda S_W W \Rightarrow S_W^{-1} S_B W = \lambda W$$

因此我們只要解出  $S_W^{-1} S_B$  的 eigen vector 即可找到投影向量

### kernel trick

這樣一樣需要使用 kernel trick，推導來源參考 kernel LDA from wiki

([https://en.wikipedia.org/wiki/Kernel\\_Fisher\\_discriminant\\_analysis](https://en.wikipedia.org/wiki/Kernel_Fisher_discriminant_analysis))

$$S_B^\phi = \sum_{j=1}^k |C_j|(m_j^\phi - m^\phi)(m_j^\phi - m^\phi)^T, m_j^\phi = \frac{\sum_{x_i \in C_j} \phi(x_i)}{|C_j|}$$

$$S_W^\phi = \sum_{j=1}^k \sum_{x_i \in C_j} (\phi(x_i) - m_j^\phi)(\phi(x_i) - m_j^\phi)^T$$

$$W = \sum_i^w \sum_j^n W_{i,j} \phi(x_j), \frac{W^T S_B^\phi W}{W^T S_W^\phi W} = \frac{\alpha^T M \alpha}{\alpha^T N \alpha}$$

$$M = \sum_{j=1}^k |C_j|(M_j - M_{all})(M_j - M_{all})^T, M_j = \frac{\sum_{x_i \in C_j} \sum_{x_j}^n K(x_i, x_j)}{|C_j|}, N = \sum_{j=1}^k K_j(I - 1_{|C_j|})K_j^T$$

這邊是實作的程式碼

```

1  def LDA(datas, labels, k, kernel=None):
2      C = np.zeros((datas.shape[0], len(np.unique(labels))))
3      for idx, j in enumerate(np.unique(labels)):
4          C[labels == j, idx] = 1;
5
6      if callable(kernel):
7          datas = kernel(datas, datas)
8
9      Mj = np.matmul(datas.T, C) / np.sum(C, axis = 0)
10     M = np.sum(datas.T, axis = 1) / datas.shape[0]
11
12     B = Mj - M[:,None]
13     SB = np.matmul(B * np.sum(C, axis = 0), B.T)
14
15     W = datas.T - np.matmul(Mj, C.T)
16     SW = np.zeros(SB.shape)
17     for group in np.unique(labels):
18         w = W[:, labels == group]
19         SW += (np.matmul(w, w.T) / w.shape[1])
20
21     try:
22         SW_inv = np.linalg.inv(SW)
23     except:
24         SW_inv = np.linalg.pinv(SW)
25
26     SBW = np.matmul(SW_inv, SB)
27     value, vector = np.linalg.eigh(SBW)
28
29     start_idx = 1 if callable(kernel) else 0
30     max_idx = np.flip(np.argsort(value))
31     W = np.concatenate([
32         vector[:, max_idx[i]][:,None]
33         for i in range(start_idx, start_idx+k)
34     ], axis=1)
35
36     return np.matmul(datas, W), W

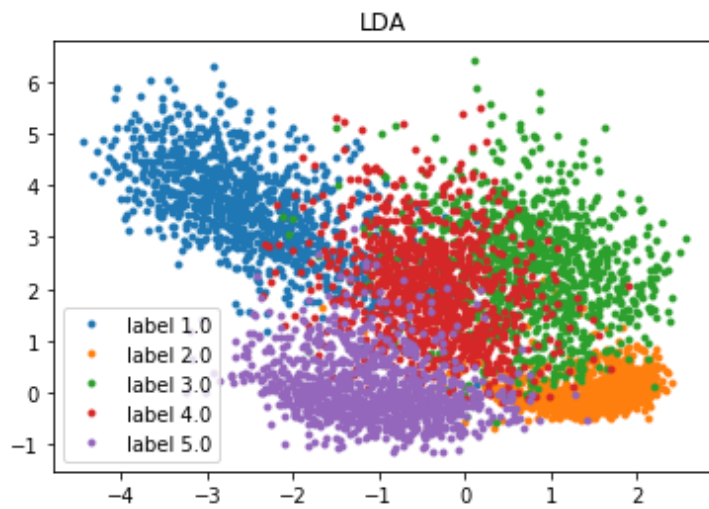
```

基本架構與 PCA 差不多，關鍵在於要找 eigen vector 的矩陣是誰  
 這樣有一個需要注意的地方是，因為需要計算矩陣的 inverse  
 因此要考量到如果不能 inverse 的時候要使用 pseudo inverse

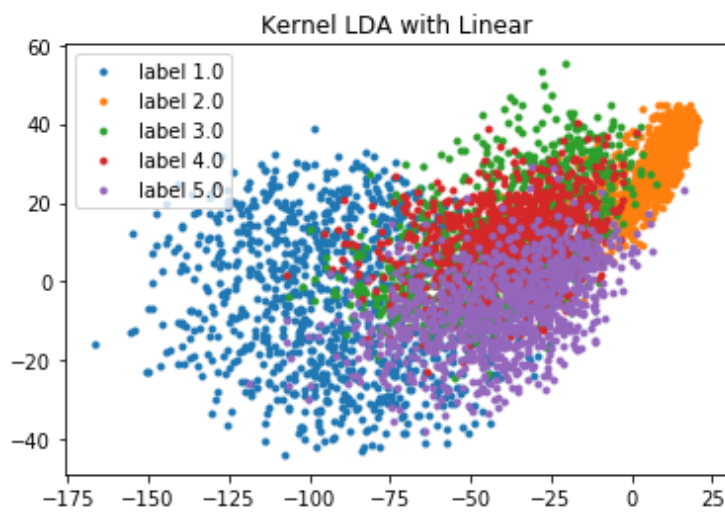


## visualization

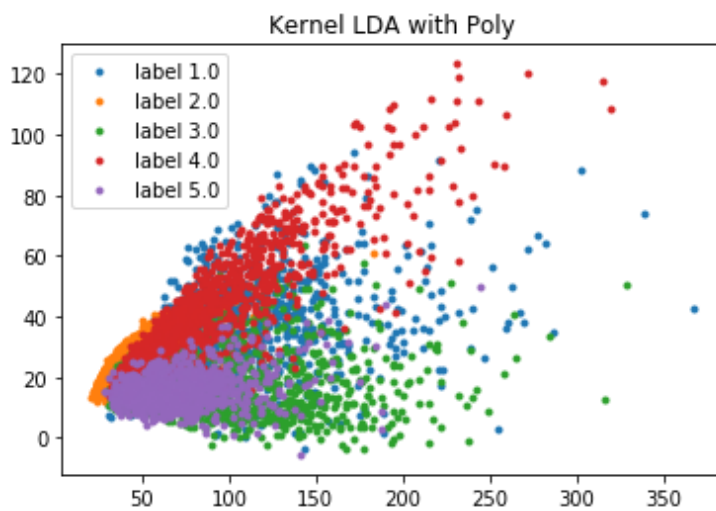
- LDA



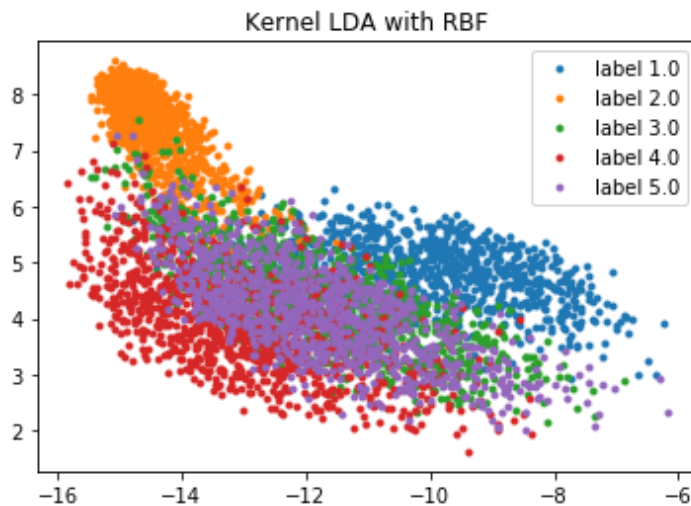
- LDA with Linear kernel



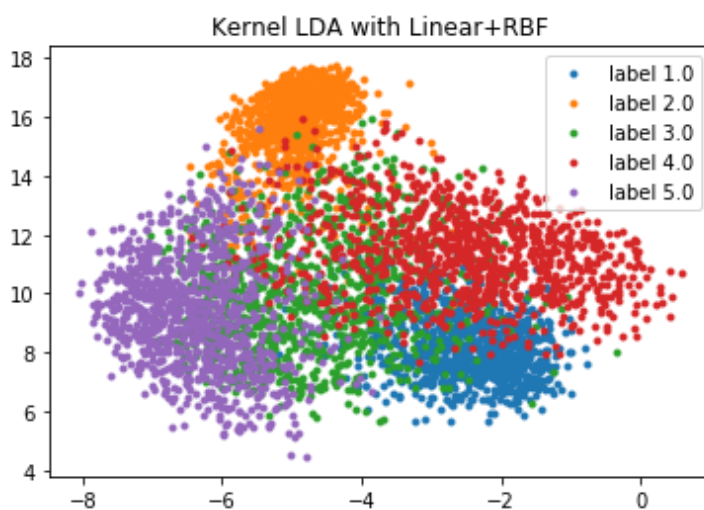
- LDA with Polynomial kernel



- LDA with RBF kernel



- LDA with Linear+RBF kernel



## my observation

在單純的 LDA (without kernel) 下，可以看到各群中心盡可能分開  
且同一群盡可能靠近，這結果確實比 PCA 好，畢竟 LDA 是知道各群的分類

加入 kernel 之後，發現到 Linear kernel 的結果並沒有跟單純的 LDA 一樣  
個人認為是我的實作有問題，因此反覆檢查與推敲也沒有找到答案

本以為 RBF 結果會不錯，但是也沒有  
反倒是 Linear+RBF 看起來比較好

## RatioCut & NormalizeCut

這邊參考我自己的 2018 ML HW6 Clustering (<https://hackmd.io/s/rJ-049ikE#steps-amp-my-implementation>)  
實作的程式，並且加入不同的 cut  
改寫成以下程式

```

1  def __normalizecut(W,D):
2      d = np.sqrt(np.linalg.inv(D))
3      return np.matmul(np.matmul(d, D-W), d)
4
5  def __generalizecut(W,D):
6      return np.matmul(np.linalg.inv(D), D-W)
7
8  __spectral_cut = {
9      0 : lambda W,D:D-W,
10     1 : lambda W,D:(D-W)/D.shape[0],
11     2 : __normalizecut,
12     3 : __generalizecut,
13 }
14 def Spectral(datas, kernel, k=None, cutType=0):
15     """
16     spectral clustering
17     parameters:
18         datas : data source n*d
19         kernel : function for compute kernel
20         k : mapping data to k-dim eigen space
21         cutType:
22             - 0 : simple cut, D-W
23             - 1 : ratio cut, (D-W)/|V|
24             - 2 : normalize cut,  $D^{-1/2}(D-W)D^{-1/2}$ 
25             - 3 : generalize cut,  $D^{-1}(D-W)$ 
26     return:
27         U : data in eigen space n*k
28         eigen values : n*1
29         eigen vector : n*n
30     """
31     if not callable(kernel):
32         raise AttributeError;
33     W = kernel(datas, datas)
34     D = np.sum(W, axis=1)*np.eye(W.shape[0])
35
36     if not cutType in __spectral_cut:
37         raise AttributeError;
38     L = __spectral_cut[cutType](W,D)
39
40     eig_values, eig_vectors = np.linalg.eigh(L)
41
42     return eig_values, eig_vectors
43
44 def GetKeigen(eigen_values, eigen_vectors, k, isMin=True, startidx=1):
45     sort_idx = np.argsort(eigen_values)
46     if not isMin:
47         sort_idx = np.flip(sort_idx)
48     return np.concatenate([
49         eigen_vectors[:,sort_idx[i]][:,None]
50         for i in range(startidx, startidx+k)
51     ], axis=1)

```

kernel 實作使用前面所寫的部份

這邊加入三種 cut 的算法，總共四種，如下

- SimpleCut
  - $\Rightarrow L = D - W$
- RatioCut
  - $\Rightarrow L_R = \frac{L}{|V|} = \frac{D-W}{|V|}$
- NormalizedCut
  - $\Rightarrow L_N = D^{-\frac{1}{2}} L D^{-\frac{1}{2}}$
- GeneralizeCut
  - $\Rightarrow L_G = D^{-1} L$

目標一樣是希望 最小化 cut，因此找出第二小之後的 k 個 eigen vectors  
為什麼不找最小的，是怕出現全連接的情況因此避開(讓所有點都成為一群，cut 就會是 0)

我將 spectral clustering 分群過程分成三步(三個 function)，方便調整與觀察

1. 計算相似度矩陣並根據 cut 找出 eigen values 和 eigen vectors
2. 根據 eigen values 挑出 k 個 eigen vectors，並組成 eigen space
3. 將 eigen space 丟入 kmeans，得到分群結果

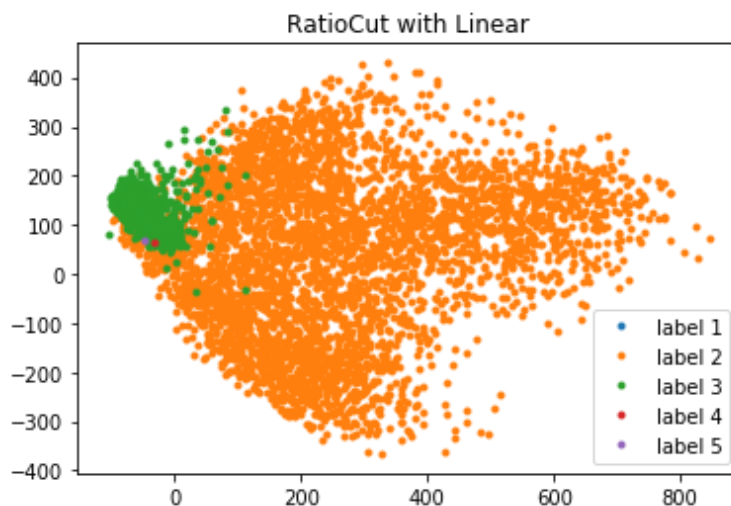
kmeans 實作直接拿 HW6 (<https://hackmd.io/s/rJ-049ikE#k-means--kernel-k-means>) 就不重複贅述

以下將跑三種不同 kernel 搭配兩種不同的 cut  
參數使用如下

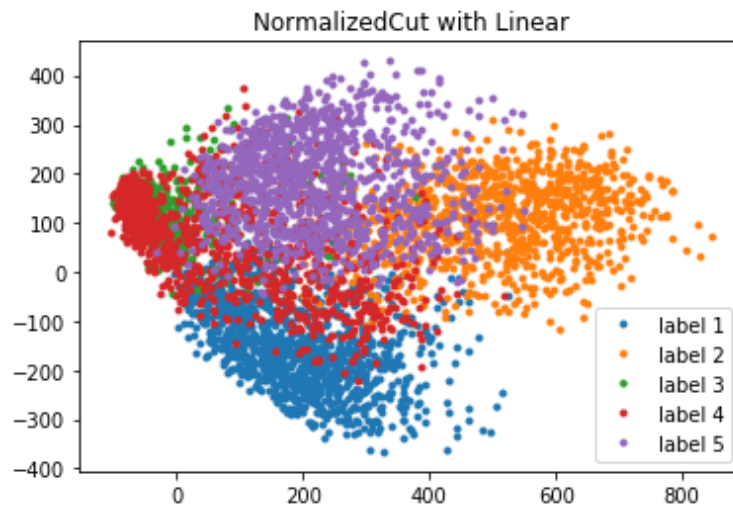
- gamma  $\gamma = 0.01$

### visualization (3 Kernel with 2 cut)

- Linear
  - RatioCut

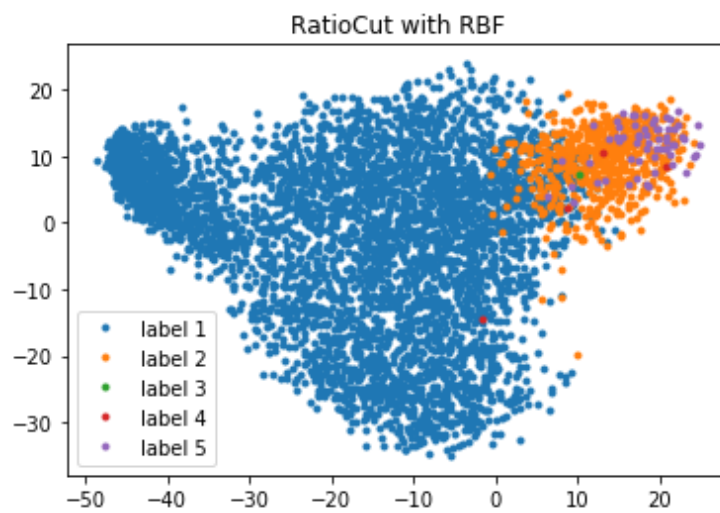


- NormalizeCut

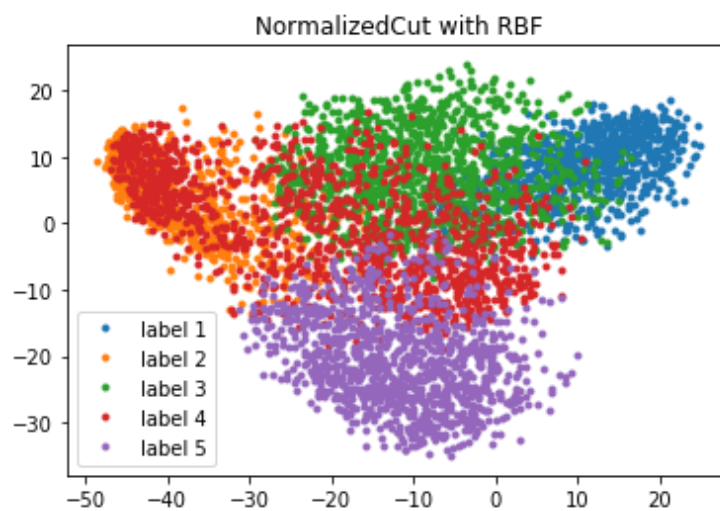


- RBF

- RatioCut

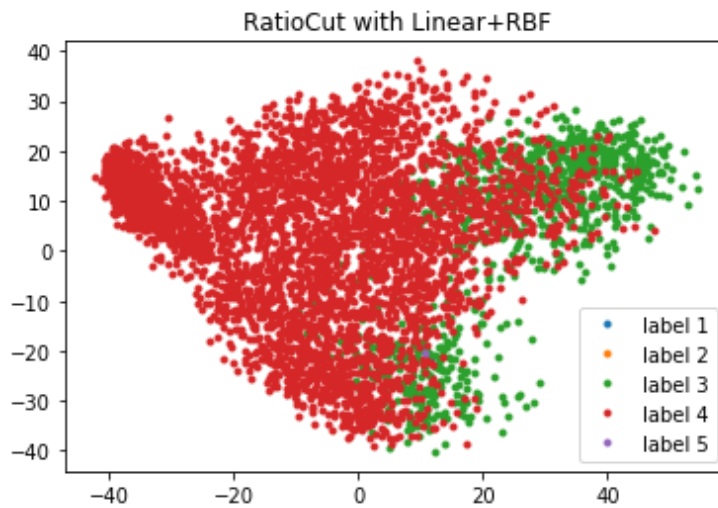


- NormalizeCut

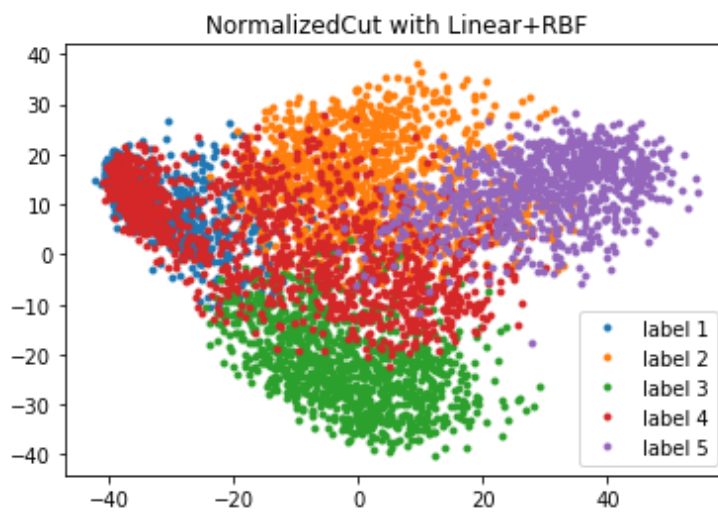


- Linear + RBF

- RatioCut



- NormalizeCut

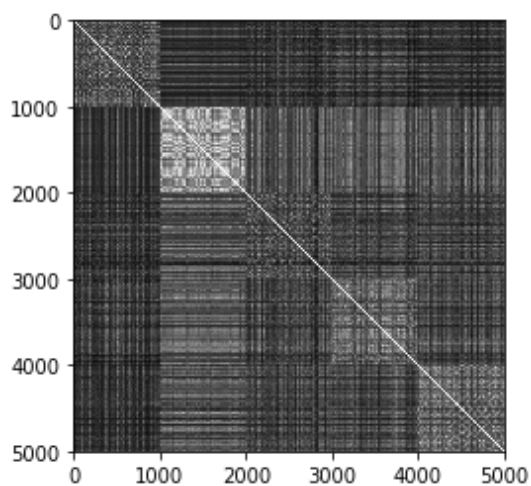


### my observation

兩種不同的 cut 的結果非常不同，NormalizeCut 能很好在圖上看到 5 個群  
而 RatioCut 則在各群之間的數量差距過大，導致直接看圖會覺得有些群消失了

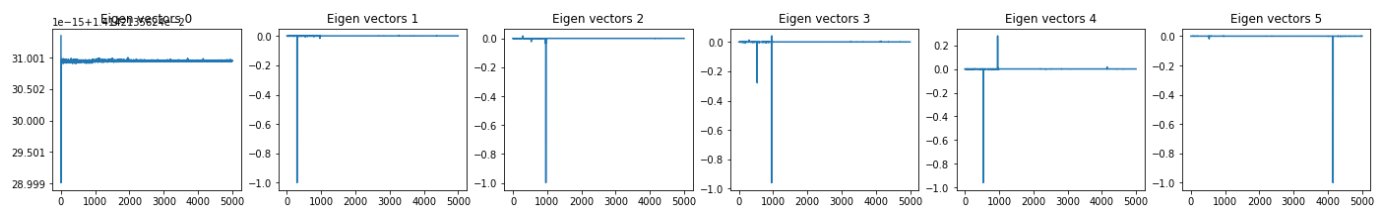
這樣的現象可以從相似度矩陣略知一二

下圖是 RBF kernel 計算得到的相似度矩陣

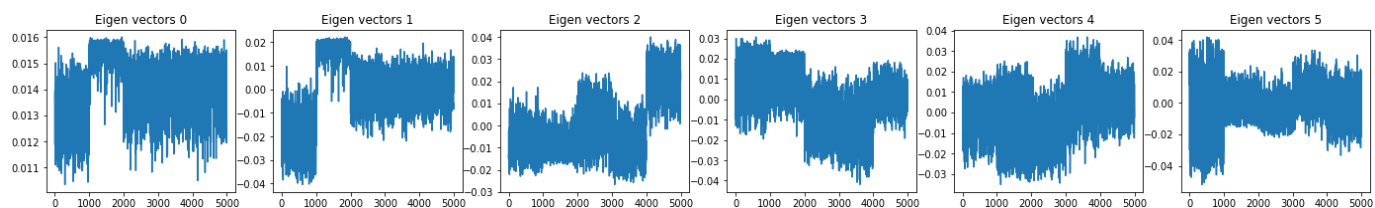




可以發現有一塊特別白(數值較高)，因此容易被分成一群  
但是在沒有 normalize 的情況下，其他非這一塊的容易直接被分在同一群



因此加入 normalize 試圖在考慮相似度下去平衡每一群的數量  
所以找到的 eigen vector 就盡量減少過度不平衡的情況



## SVC result visualization

直接引用 HW5 的程式碼，將分群得出的結果與 support vector 呈現出來

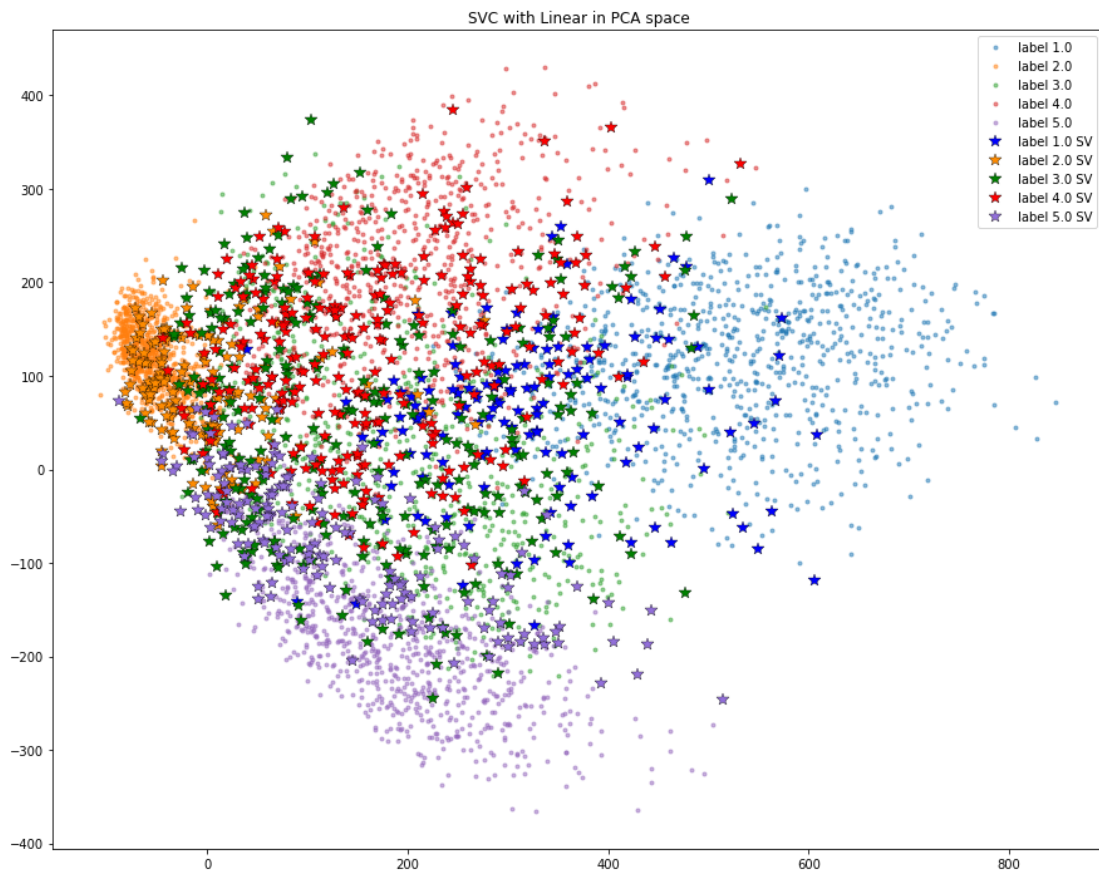
```
1  import sys
2  sys.path.append("../HW5")
3  #import libsvm.python.svm as svm
4  import libsvm.python.svmutil as svmutil
5
6  def sparse_matrix(x, iskernel=False):
7      row = x.shape[0]
8      col = x.shape[1]
9      idx_offset = 1
10     if (iskernel):
11         x = np.append(np.linspace(1,row,row), x).reshape(col+1,row).T
12         idx_offset = 0
13     x = [{idx+idx_offset:x[i][idx] \
14           for _,idx in np.ndenumerate(np.argwhere(x[i]!=0))} \
15           for i in range(x.shape[0])]
16     return x
17
18     __svc_probs = {}
19     __svc_models = {}
20
21     # establish svm problems
22     for kernelname, kernel in __kernels.items():
23         if not kernelname in __svc_probs:
24             __svc_probs[kernelname] = svmutil.svm_problem(
25                 train_label,
26                 sparse_matrix(
27                     __kernels[kernelname](train_data, train_data),
28                     iskernel=True),
29                 isKernel=True)
30
31     # train model with problem
32     __svc_c = {
33         'Linear' : 0.01,
34         'Poly' : 1,
35         'RBF' : 100000,
36         'Linear+RBF' : 1000
37     }
38     for kernelname, prob in __svc_probs.items():
39         if not kernelname in __svc_models:
40             __svc_models[kernelname] = svmutil.svm_train(
41                 prob,
42                 '-t 4 -c {}'.format(__svc_c[kernelname]))
```

全部一律利用 precomputed kernel data 的方式，已保證 kernel 參數與 PCA 與 LDA 使用的 kernel 參數一致，之後取得 model 當中的 `sv_indices` 就可以知道是哪些資料點成為了 support vector，至於 C-SVC 的參數 `c` 則是參考 HW5 的，盡可能好的參數

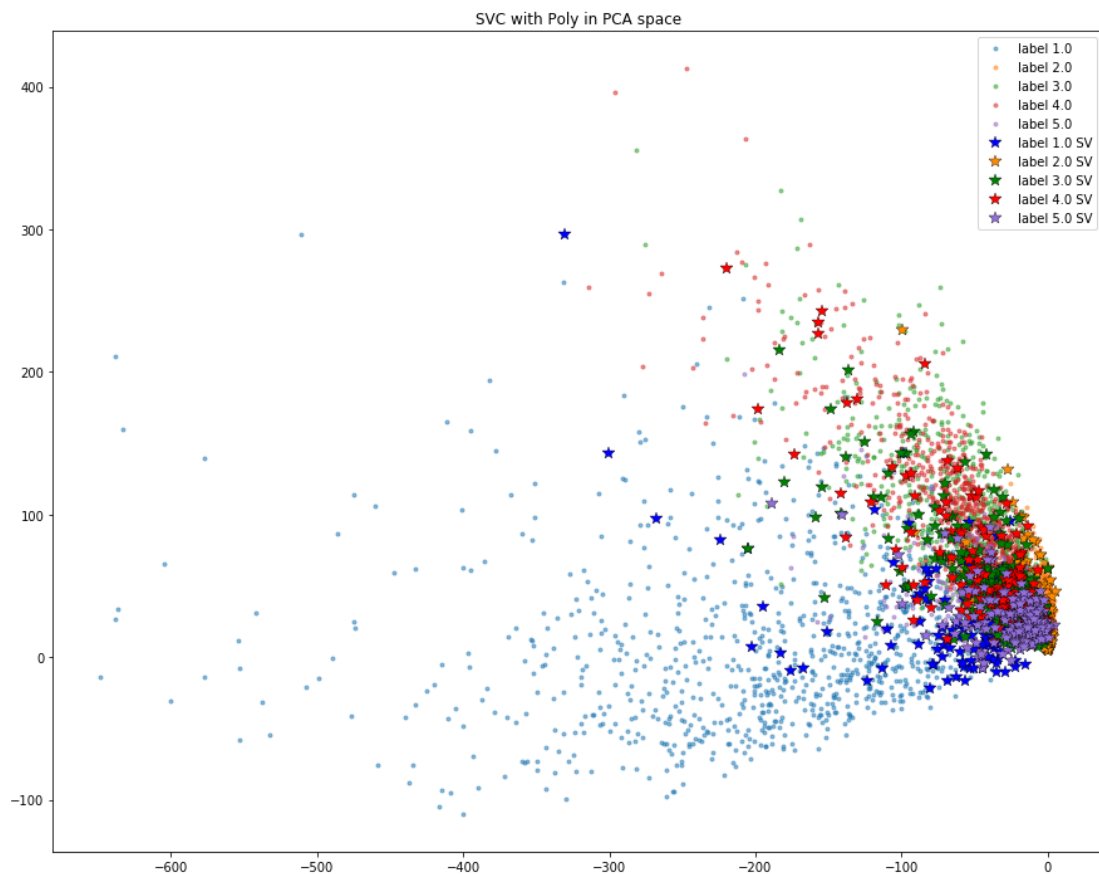


## In PCA space

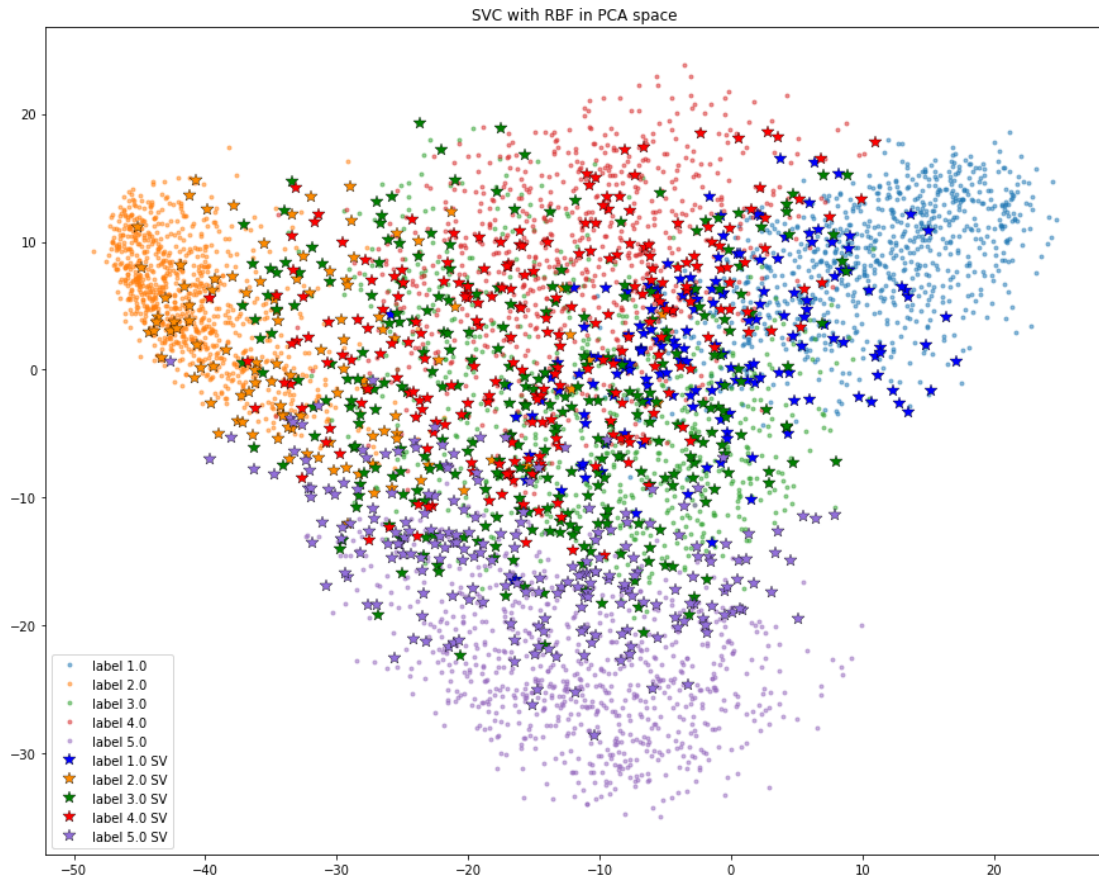
- Linear



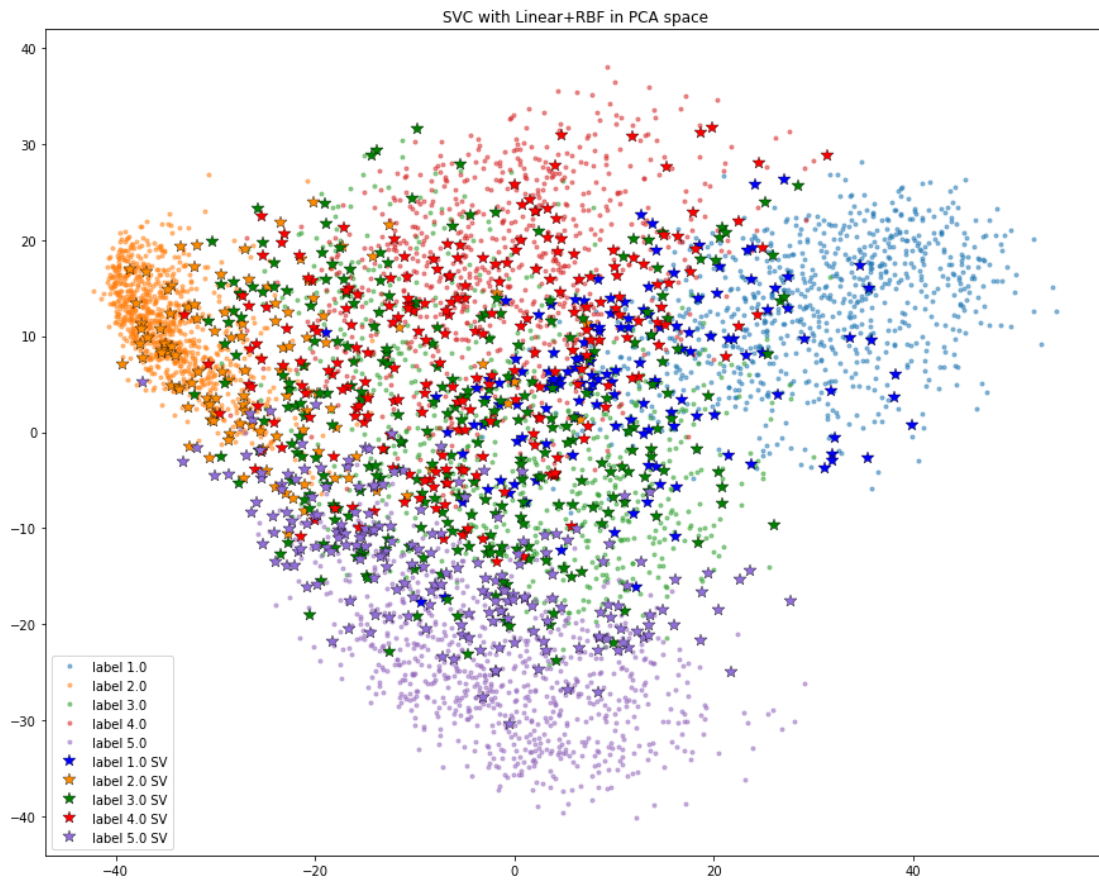
- Polynomial



- RBF

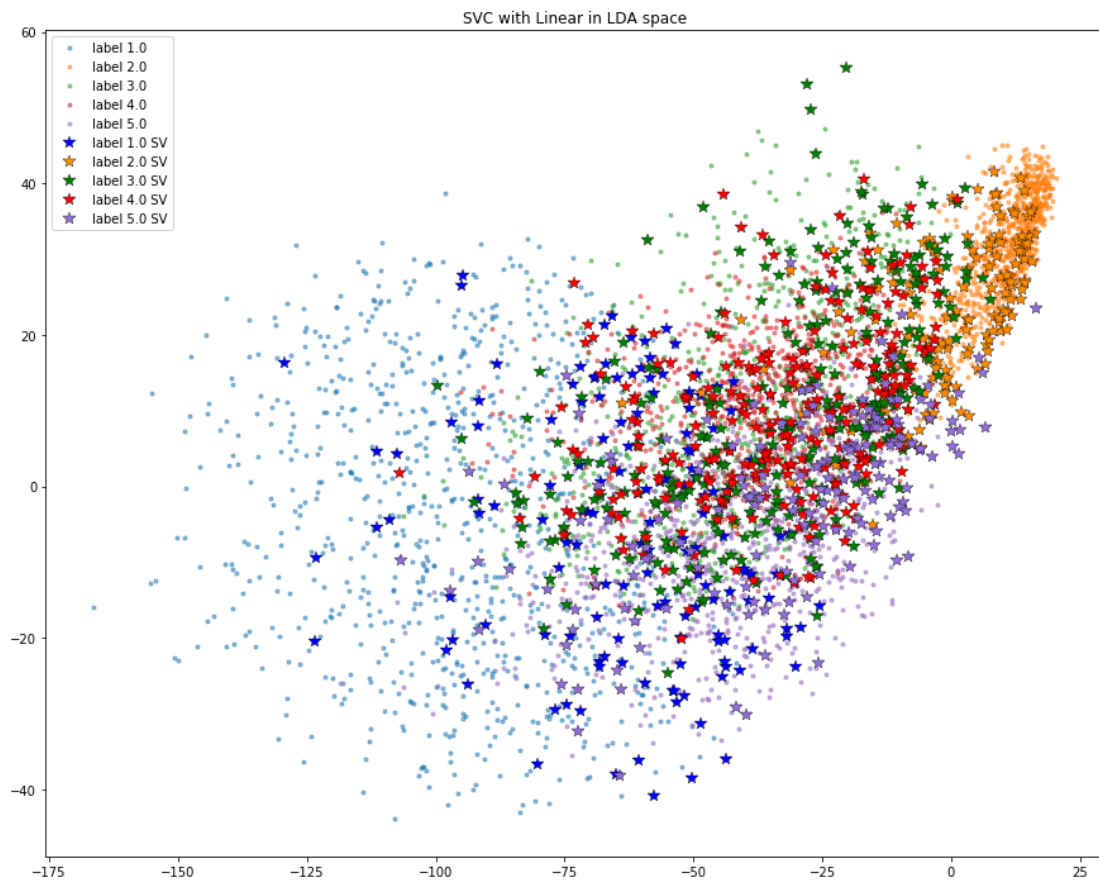


- Linear+RBF

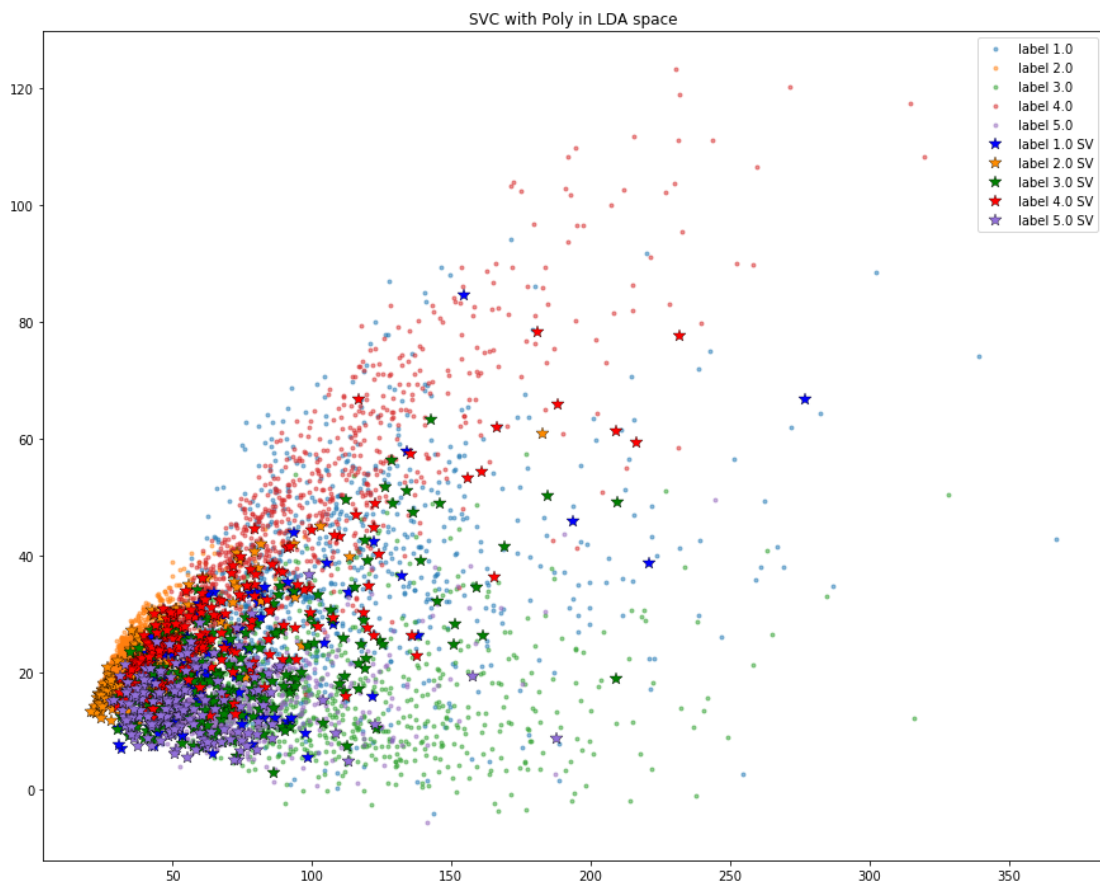


## In LDA space

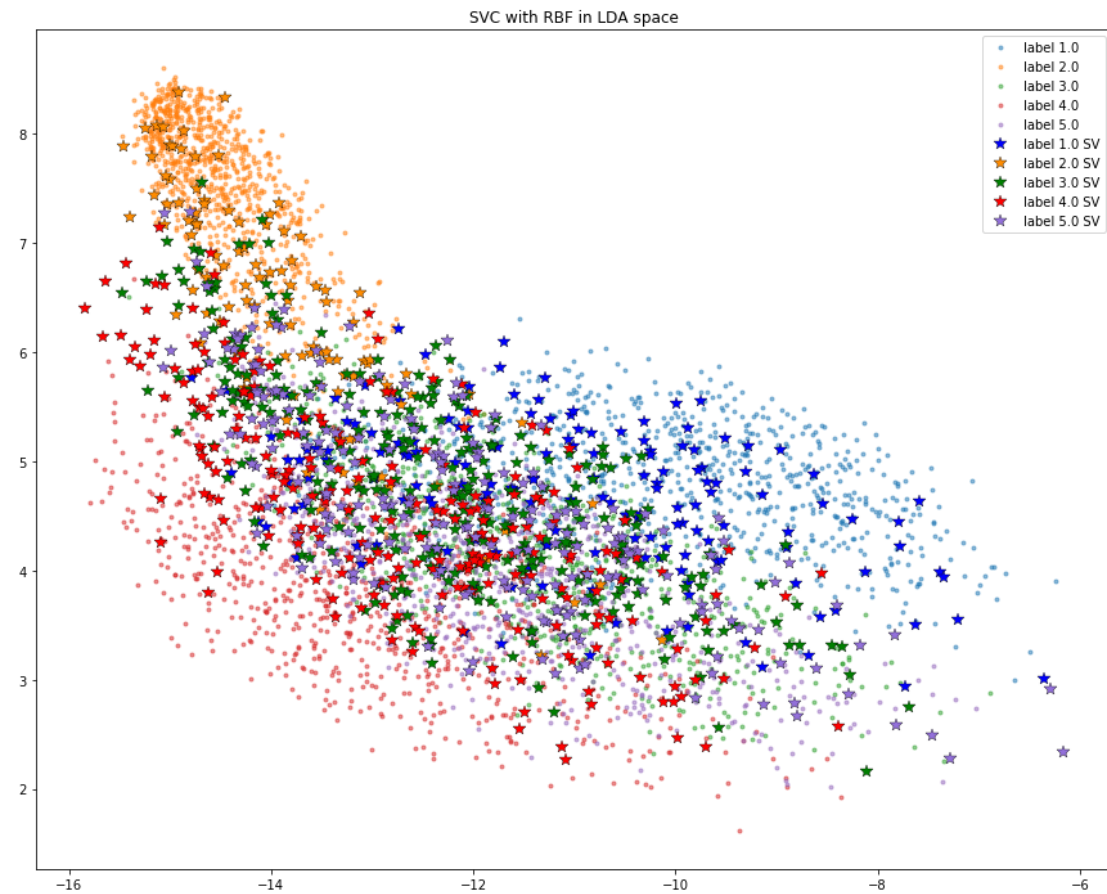
- Linear



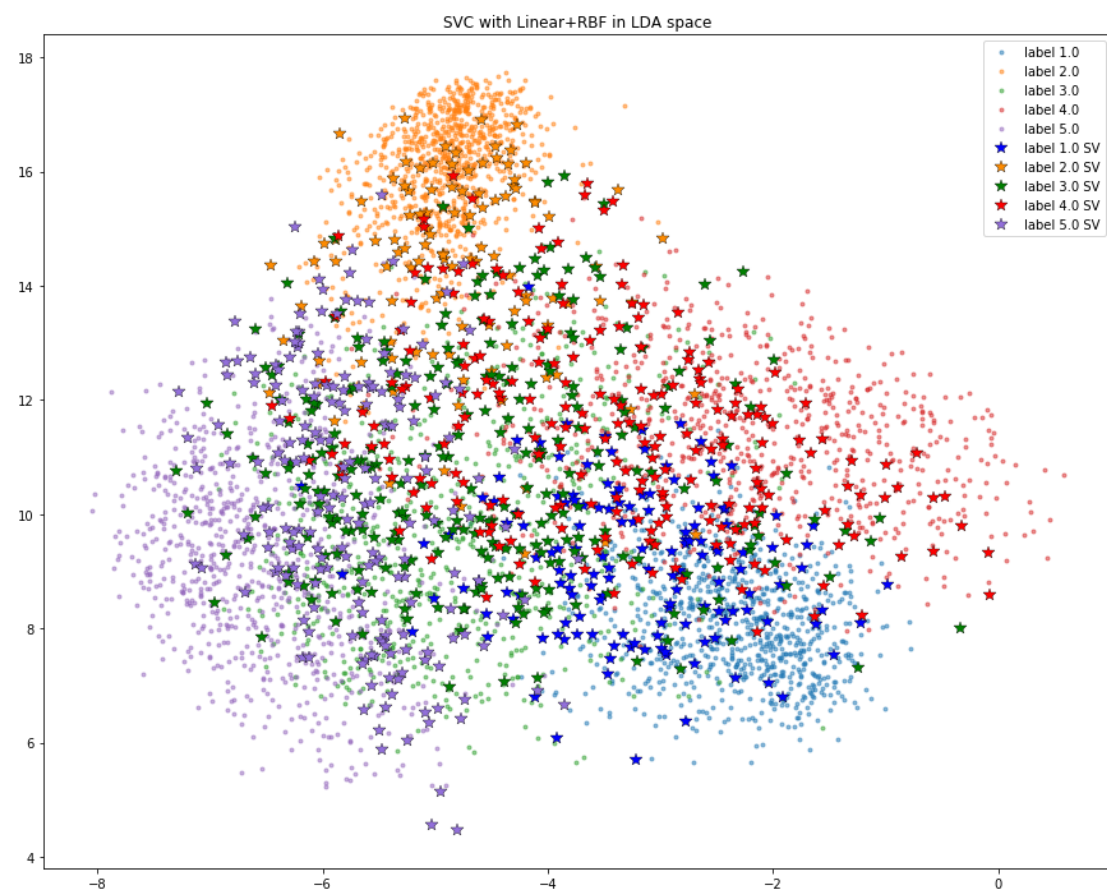
- Polynomial



- RBF



- Linear+RBF





## my observation

這邊最大的成就感來自於我可以看到 SVC 分出來的樣子  
不然在 HW5 就是只能看到 accuracy，卻無法有視覺上去呈現分的好與否  
藉由 PCA/LDA 的降維，可以將資料視覺化

可以看到 support vector 有在當群的邊界，常出現在群與群的分界點  
當然也有一些很神奇的位置會出現 support vector

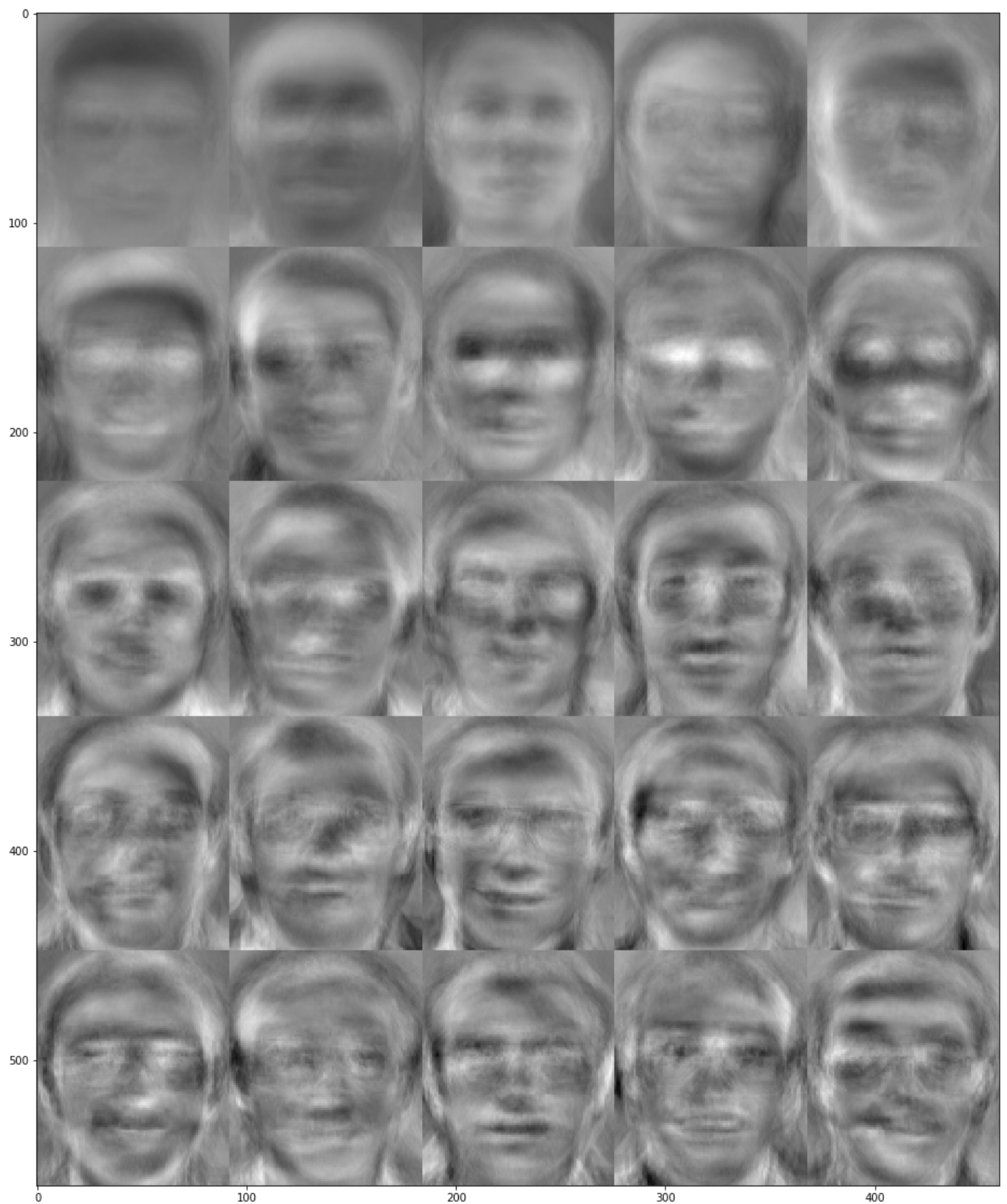
## eigenfaces from att\_faces

先將人臉圖像讀進來存成 nparray 方便後續處理

```
1  import PIL.Image
2  att_faces_data = []
3  for i in range(1,41):
4      for j in range(1,11):
5          with PIL.Image.open(
6              './att_faces/s{}/{}.pgm'.format(i, j)
7          ) as im:
8              att_faces_data += [np.array(im).reshape(1,-1)]
9
10 att_faces_data = np.concatenate(att_faces_data, axis = 0)
```

att\_faces\_data 為 400x10304 的矩陣 (400 張人臉，一張 112x92)  
使用 PCA 找到 25 個主成份也就 eigenface (特徵臉)

```
11 face_pcaspace, face_eigen, face_eva, face_evc = PCA(att_faces_data, 25)
12 showFaces(face_eigen.T, col=5)
```

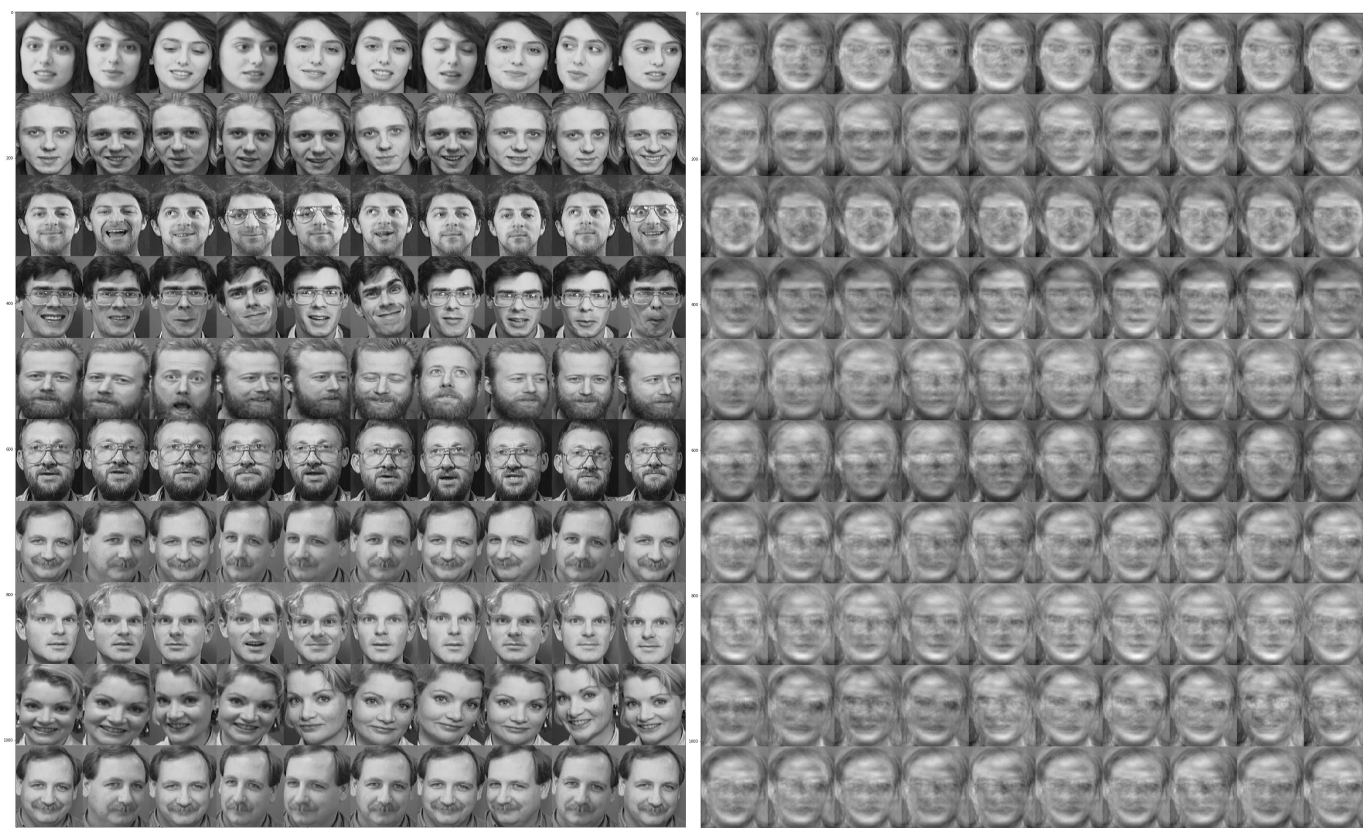


並且將原本的人臉投影到特徵臉在投影回來  $\text{face} * \text{eigenfaces} * \text{eigenfaces}^T$   
只隨機挑選 10 位人臉的不同樣子 (表情，眼鏡穿戴與否，臉轉動等)

```

13 # random pick subject
14 # random pick subject
15 random_faces = np.concatenate([
16     att_faces_data[sj*10:(sj+1)*10]
17     for sj in np.random.randint(0,40, size=10)]
18     , axis=0)
19
20 showFaces(random_faces)
21 showFaces(np.matmul(
22     np.matmul(random_faces, face_eigen),
23     face_eigen.T)
24 )

```



左邊為原本的人臉，右邊投影到 eigenfaces 上再重建回來

### my observation

結果還蠻有趣的，如果觀察 eigenfaces 會發現有人臉的輪廓

同時也感覺到有不同五官出現，只是較為模糊

找到這樣特徵臉可以幫助我們做人臉辨識，可以將 112x92 共 10340 維度的資料弄到只有 25 維，所代表含意是在 25 張特徵臉成份多寡

但把 400 張人臉 (40 個人) 投影到 25 張特徵臉是否會導致每個人的獨特特徵完全消失導致只能辨識這是人臉但無法辨識這是誰的臉？

這邊可以觀察重構後的人臉(投影到特徵臉空間在投影回原本的資料空間)

可以發現右邊重構的人臉還是能夠看出是不同人，只是清晰的五官都模糊了

代表在特徵臉空間的 25 維向量是可以拿來做分類並達到 身份辨識

而且除此之外，也有一定的容錯空間，每個人有十張不同姿態的人臉  
有沒表情，有表情，有眼鏡沒眼鏡，傾斜等等  
但可以看到重構後同一個人的十張臉幾乎都長的差不多  
這也代表在特徵臉空間同一個人不同姿態的向量不會距離太遠

藉由這樣的例子也可以感受到 PCA 的厲害  
找出該資料的主要成分，並藉由該成份去表達原本資料  
同時去除無用資料與雜訊只關注我們想關注的主要特徵