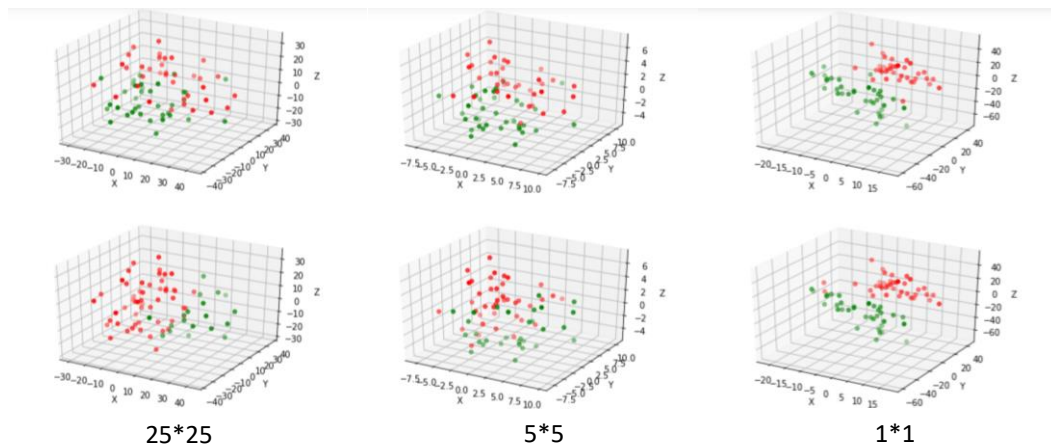


Lab1 PCA+K-MEANS

一、壓縮和讀檔

```
38 #壓縮圖片
39 def compressImage(srcPath,dstPath):
40     for filename in os.listdir(srcPath):
41         if not os.path.exists(dstPath):
42             os.makedirs(dstPath)
43
44         #完整的檔案或資料夾路徑
45         srcFile=os.path.join(srcPath,filename)
46         dstFile=os.path.join(dstPath,filename)
47
48         #如果是檔案就處理
49         if os.path.isfile(srcFile):
50             if os.path.isfile(dstFile):
51                 try:
52                     img=Image.open(srcFile)
53                     dimg=img.resize((1,1),Image.ANTIALIAS) # 設定壓縮尺寸
54                     dimg.save(dstFile)
55                 except Exception:
56                     print(dstFile+"fail")
57
58         #如果是資料夾就遞迴
59         if os.path.isdir(srcFile):
60             compressImage(srcFile, dstFile)
61 ~
```

一開始在做 PCA 時因為 memory 不夠，無法儲存這麼大的投影矩陣，所以決定在讀檔處理圖片時，就先壓縮圖片。



嘗試將不一樣大的圖片統一壓縮成不同的大小後，發現壓縮得越小，在分類的表現也會越好，同時注意到當我把圖片壓縮成 1*1 時，經過 PCA 降維處理後，在三維空間的分布跟前面相比呈現比較大的區別。

```
91 compressImage("./sum/green/", "./sum/green_compression/") # 壓縮綠色芒果圖片
92 path= './sum/green_compression/'
93 files=os.listdir(path)
94 img_list = []
95
96 #讀檔
97 for file in files:
98     p=path+file
99     img = cv2.imread(p)
100     img_2D = np.reshape(img,(-1,3))
101     img_1D = np.reshape(img,(-1,1)) # 將原先三維的數值展開成一維
102     img_1D = img_1D.flatten()
103     img_list.append(img_1D)
104 ~
```

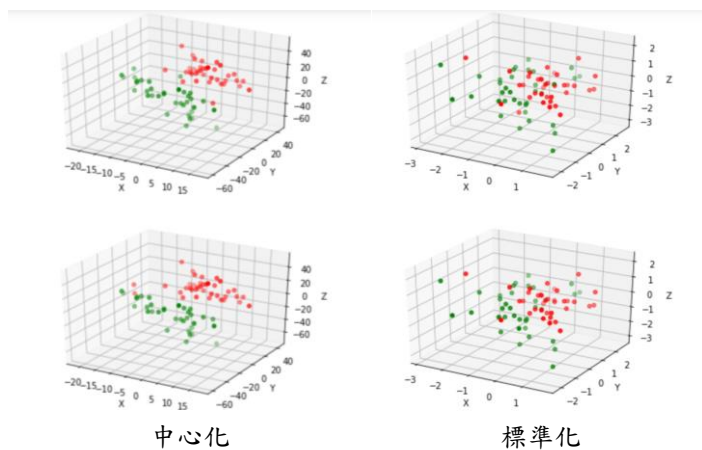
原本我是將每張圖片降維成二維 array，大小是一張圖片壓縮過後的 element 個數*每個樣本的屬性值個數(3)，然後對每張照片分別進行 PCA 處理，但結果顯示紅綠分布很明顯混雜在一起，所以後來改成先將每一張照片樣本的所有元素拉為一為陣列，再將代表每一張照片的 array 組成 matrix，最後一起丟入 PCA 進行處理。

二、PCA

```
11 #標準化
12 def standardization(data):
13     mu = np.mean(data, axis=0)
14     sigma = np.std(data, axis=0)
15     return (data - mu) / sigma

17 #pca
18 def pca(XMat, k):
19     m, n = np.shape(XMat)
20     average = np.mean(XMat, axis=0)
21     #data_adjust = XMat - average
22     data_adjust = standardization(XMat)
23     #print(data_adjust.shape)
24     covX = np.cov(data_adjust.T)
25     #得到協方差矩陣
26     eigenValue, eigenVec = np.linalg.eig(covX)
27     #取特徵值和特徵向量
28
29     if k > n:
30         print('k must be lower than feature number')
31         return
32     else:
33         selectVec = eigenVec[:, :k]
34         #選取K個特徵向量，建立投影矩陣
35         finalData = np.dot(data_adjust, selectVec)
36         #投影
37         finalData = standardization(finalData)
38         #reconData = np.dot(finalData, selectVec.T) + average
39         return finalData, XMat
```

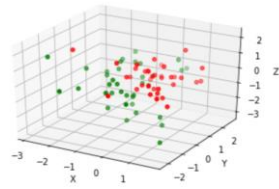
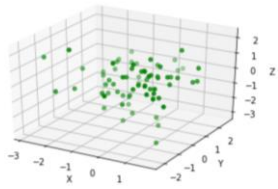
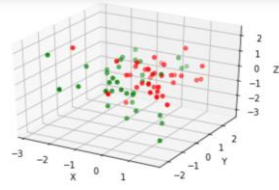
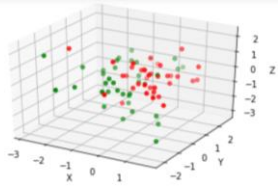
針對要準備進行PCA處理的數據，一開始我是採中心化處理，讓他們與平均相減，後來我又嘗試標準化的處理，分類效果也有提升。



三、K-MEANS

```
66 def kmeans_xufive(ds, k):
67     m, n = ds.shape
68     result = np.empty(m, dtype=np.int)
69     cores = np.empty((k, n))
70     cores = ds[np.random.choice(np.arange(m), k, replace=False)]
71     #print(cores)
72     #cores = [
73     #    [np.random.randint(-3,1), np.random.randint(-2,2), np.random.randint(-3,2)]
74     #    for i in range(2)]
75     #]
76     #隨機生成質心
77
78     while True: #迭代計算
79         d = np.square(np.repeat(ds, k, axis=0).reshape(m, k, n) - cores)
80         distance = np.sqrt(np.sum(d, axis=2))
81         index_min = np.argmin(distance, axis=1)
82         #計算每個樣本距離k個質心的距離
83         #標記每個樣本距離最近的質心
84
85         if (index_min == result).all():
86             return result, cores
87         #如果樣本分類沒有改變
88         #則回歸分類結果和質心資料
89
90         result[:] = index_min
91         #重新分類
92         for i in range(k):
93             items = ds[result==i]
94             #找出對應當前質心的分類樣本
95             if len(items) != 0:
96                 cores[i] = np.mean(items, axis=0)
97             #取分類樣本座標平均值作為新質心的位置
```

原先我是隨機生成兩個點做為分類的起始點，但後來發現這樣很容易讓分類不均，甚至遇到全部都分成同一類，所以後來我改成隨機取 72 個樣本數中的兩個點做為起始點，分類效果明顯好了很多。



隨機在空間中取點

隨機從樣本點取點