

3 Task 3

The code of principal component analysis:

```
1 def zeroMean(dataMat):
2     meanVal=np.mean(dataMat, axis=0)
3     newData=dataMat-meanVal
4     return newData,meanVal
5
6 def pca(dataMat,n):
7     newData,meanVal=zeroMean(dataMat)
8     covMat=np.cov(newData, rowvar=0)
9
10    eigvals,eigVects=np.linalg.eig(np.mat(covMat))
11    eigValIndice=np.argsort(eigvals)
12    # print(eigValIndice)
13    n_eigvalIndice=eigValIndice[ -1:-(n+1):-1] # pick n values
        from last to start
14    # print(n_eigvalIndice)
15
16    # choose the smallest one
17    # n_eigvalIndice=1 # 1 is the index of the smallest
        eigenvalue
18
19    n_eigvect=eigVects[ :,n_eigvalIndice]
20    # print(n_eigvect)
21    lowDDataMat=newData*n_eigvect
22    reconMat=( lowDDataMat*n_eigvect.T)+meanVal
23    return eigvals,eigVects,lowDDataMat ,reconMat
```

3.1 a)

The result of eigenvalue decomposition of the corresponding sample covariance matrix is [1.5, 0.16666667]

3.2 b

1D representation under the **largest** eigenvalue is `[[-1.06066017],[-1.06066017],[1.06066017],[1.06066017]]`. And 1D representation under the **smallest** eigenvalue is `[[-0.35355339],[0.35355339],[-0.35355339],[0.35355339]]`

3.3 c

The code of disance computation

```

1 def get_dis_pairs(M):
2     lens=M.shape[0]
3     result=np.zeros((lens,lens))
4     for i in range(lens):
5         j=i+1
6         while(j<lens):
7             result[i][j]=np.linalg.norm(M[i]-M[j])
8             j=j+1
9     return result

```

After euclidean distance computation, we can see the result from Fig.2. The distance between points in original matrix and reconstructed matrix after PCA are the same. But, the distance between points after PCA changed due to the dimension decreasing.

[[0. 0.70710678 2.12132034 2.23606798]	O_Matrix
[0. 0. 2.23606798 2.12132034]	
[0. 0. 0. 0.70710678]	
[0. 0. 0. 0.]]	
[[0. 0. 2.12132034 2.12132034]	1D
[0. 0. 2.12132034 2.12132034]	
[0. 0. 0. 0.]]	
[0. 0. 0. 0.]]	
[[0. 0.70710678 2.12132034 2.23606798]	2D
[0. 0. 2.23606798 2.12132034]	
[0. 0. 0. 0.70710678]	
[0. 0. 0. 0.]]	

Figure 2: Euclidean distance

3.4 d)

In new dataset. the eigenvalue decomposition of the corresponding sample covariance matrix is $[1.5, 0.16666667]$, the same as in the old dataset. when converting two datasets into 1-D form, points in their own graph are lies on the same places in Fig.3 they represent same information after principal component analysis. These two datasets represent different information on the original data space, but show the same data characteristics in one dimensional space after PCA process.

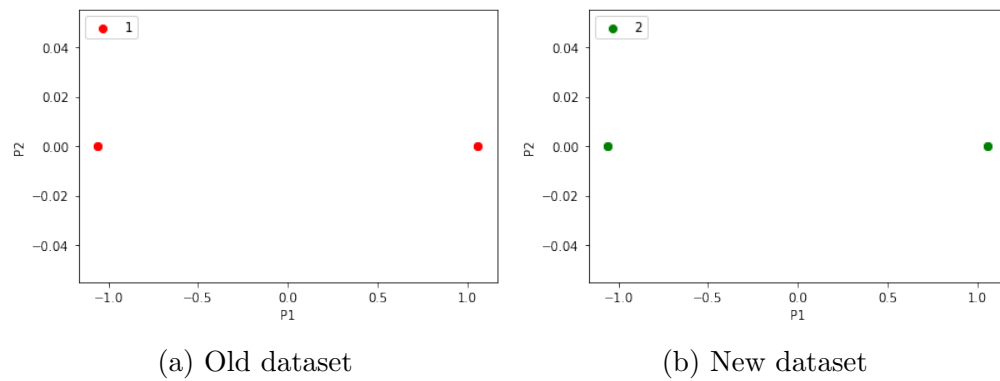


Figure 3: Points in 2 datasets