

Homework Set 4

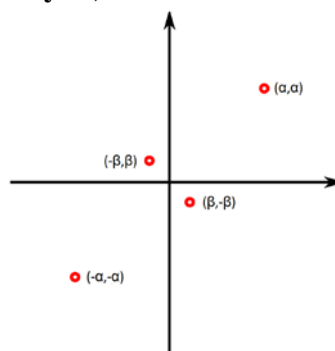
Problem 1 (Principal component analysis)

Joe's first job as a statistician was to analyze the countries that won medals in the 2012 London Olympic Games. He started his task by collecting data from 84 countries that had won at least one medal. The data's features (dimensions) for each country (sample) consisted of the number of medals won, the population size (in hundred thousands) and the GDP (in millions of US dollars) all taken in 2012. To start his analysis, Joe ran PCA on the data. The first eigenvector explained 99.9% of the data's variance. Joe's boss ran PCA by himself and he found that the first eigenvector explained only 74.0% of the data's variance. They compared what they both had done but the only difference they found was that Joe's boss used the GDP values in billions of dollars while Joe's GDP values were in millions of dollars.

Answer the following questions:

- Why did they obtain different results?
- How should Joe and his Boss have treated the data before the PCA?

Problem 2 (Principal component analysis)



- Compute the covariance matrix $\text{cov}(\mathbf{X})$ of the dataset composed of the four points shown in the figure $[(\alpha, \alpha), (-\alpha, -\alpha), (-\beta, \beta), (\beta, -\beta)]$.
- Compute the eigenvectors and eigenvalues of $\text{cov}(\mathbf{X})$. Discuss how the eigenvectors and eigenvalues are affected by α and β .

Problem 3 (Principal component analysis)

PCA decomposition is a common approach for denoising. Suppose a signal \mathbf{X} is measured. Assume the signal is additive. It can be decomposed into 2 matrices by using PCA, i.e., $\mathbf{X} = \mathbf{X}_S + \mathbf{X}_N = \mathbf{E}_S \mathbf{\Sigma}_S (\mathbf{F}_S)^T + \mathbf{E}_N \mathbf{\Sigma}_N (\mathbf{F}_N)^T$, where \mathbf{X}_S and \mathbf{X}_N respectively contain "information" and "noise" of the data. Typically noise is associated with a low level of variance in data. Hence, in the PCA decomposition, the first n eigenvectors that are associated with high variance are usually assumed to be the information. In contrast, the rest eigenvectors are assumed to be the noise. If one can identify the correct intrinsic dimensionality (e.g., correct n) of the data, he/she should expect to be able to filter out the noise while still capturing the important information in the data by reconstructing the signal \mathbf{X}_S using only the first n eigenvectors.

This problem helps you practice the denoising using PCA. A dataset "04HW2_noisy.mat" consists of 1,965 gray-level images. These images are 20×28 in dimension and are distorted by adding Gaussian noises to each pixel with a standard deviation of 25. Import the dataset into MATLAB. To view the 10th image, type:

```
>>colormap gray  
>>imagesc(reshape(X(:, 10), 20, 28))'
```

Finish the following tasks:

- a) Apply PCA to the noisy data. Suppose the correct intrinsic dimensionality of the data is 10 (i.e., only the first 10 eigenvalues and eigenvectors contain information; the rest eigenvalues and eigenvectors contain noise). Compute reconstructed images using the top 10 eigenvectors and plot the 10th, 121st, 225th, 318th, and 426th original and reconstructed images.
- b) Repeat part a), assuming the correct intrinsic dimensionality of the data is 2 and 30.
- c) Determine the best intrinsic dimensionality of the dataset using the techniques learned from the class. You will need to check all the possible dimensionality. Explain the approaches you use and the reasons of choosing the best intrinsic dimensionality in details.