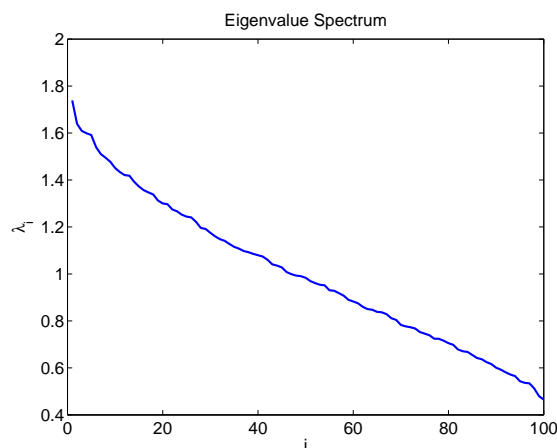


## Series 2, March 9-10, 2017 (Principal Component Analysis)

### Problem 1 (PCA Theory):

Given a dataset  $\mathbf{X} \in \mathbb{R}^{D \times N}$  (observations as columns), where  $D$  is the number of dimensions and  $N$  is the number of observations, a linear transformation using an orthonormal matrix is applied to make a *change of basis*, to obtain a (usually) lower-dimensional representation of the dataset denoted by  $\tilde{\mathbf{Z}} \in \mathbb{R}^{K \times N}$ .  $\tilde{\mathbf{Z}}$  together with the basis (and the shift), is then used to reconstruct a compressed version of the data. This can for example be applied to compress images or visualize high-dimensional data by projecting to a lower dimensional space.

1. We begin by reviewing the steps of applying PCA to a dataset  $\mathbf{X}$ . Please complete each step below by providing the appropriate formula to compute the desired quantity.
  - (a) Define the zero-mean dataset  $\bar{\mathbf{X}}$  in terms of the original dataset  $\mathbf{X}$ . For this purpose, use  $\mathbf{M} = [\bar{\mathbf{x}}, \dots, \bar{\mathbf{x}}] \in \mathbb{R}^{D \times N}$  where  $\bar{\mathbf{x}}$  is the sample mean.
  - (b) Define the covariance matrix  $\Sigma$  in terms of the zero-mean dataset  $\bar{\mathbf{X}}$ .
  - (c) Write down the eigen decomposition of the covariance matrix  $\Sigma$ , in terms of the eigenvector matrix  $\mathbf{U}$  (eigenvectors as columns) and the diagonal matrix of eigenvalues  $\Lambda$ .
  - (d) Define the dataset in the new basis  $\tilde{\mathbf{Z}}$  via a transformation of  $\bar{\mathbf{X}}$ . (Note: Assume we want to keep only the  $K$  dimensions of the transformed dataset and that the eigenvectors in  $\mathbf{U}$  have already been sorted according to the corresponding eigenvalues, in decreasing order.)
  - (e) How can the data (approximation)  $\tilde{\mathbf{X}}$  be reconstructed?
  - (f) Prove that the (squared) reconstruction error is the sum of the lowest  $D - K$  eigenvalues of the covariance matrix.
2. Assume we have applied PCA to some dataset ( $D = 100$ ). We observe the following eigenvalue spectrum of the covariance matrix of the data. ( $\lambda_i$ : eigenvalues)



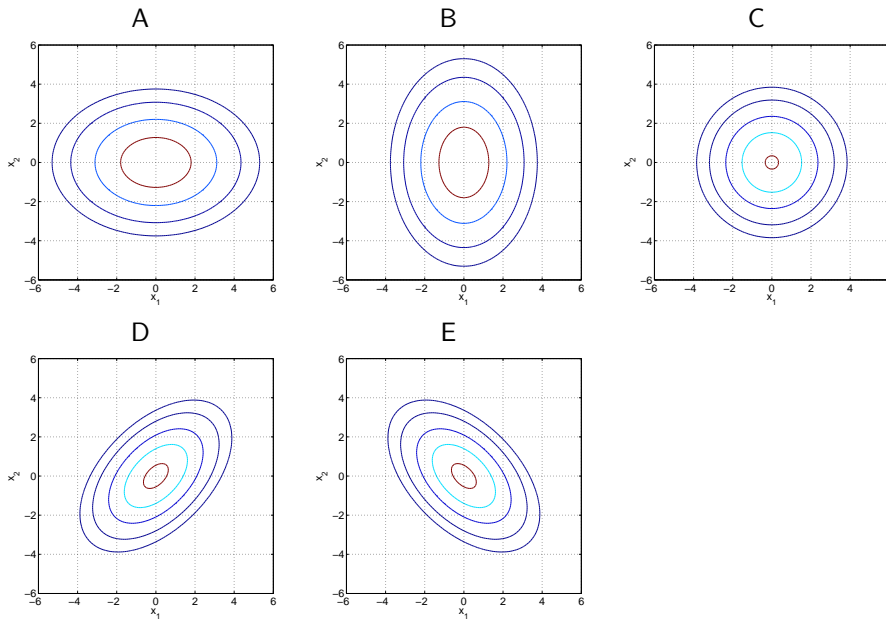
- (a) Is the intrinsic dimensionality of this dataset low or high? Why?
- (b) Can this dataset be expressed in few dimensions with low approximation error? Why?
- (c) If yes, which dimensionality (approximately) should be chosen for the transformed dataset and why?

3. Assume you have observed 2D data  $\mathbf{X} \in \mathbb{R}^{2 \times N}$  (observations as columns). The first row of  $\mathbf{X}$  corresponds to the first dimension  $x_1$ , the second row corresponds to  $x_2$ . For each of the three covariance matrices  $\mathbf{C}_{\mathbf{X}}$  below, please choose the iso-line plot (A-E) corresponding to the covariance matrix. (Note the axis labels on the figures)

1.  $\begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$  Answer: ( )

2.  $\begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$  Answer: ( )

3.  $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  Answer: ( )



4. PCA transforms a dataset  $\mathbf{X}$  into a dataset  $\mathbf{Z} = \mathbf{A}^\top \mathbf{X}$  by defining a new basis using the eigenvectors of the covariance matrix  $\Sigma_{\mathbf{X}}$  of the dataset  $\mathbf{X}$ . With this particular choice of a new basis, the covariance matrix  $\Sigma_{\mathbf{Z}}$  of the transformed dataset  $\mathbf{Z}$  is *diagonalized*.
- Please explain in words, why we desire the covariance matrix of the transformed dataset to be diagonal.
  - Show that  $\Sigma_{\mathbf{Z}} = \mathbf{A}^\top \Sigma_{\mathbf{X}} \mathbf{A}$ , i.e., that the covariance matrix  $\Sigma_{\mathbf{Z}}$  of the transformed dataset can be written in terms of the covariance matrix  $\Sigma_{\mathbf{X}}$  of the original dataset.
  - Show that the choice  $\mathbf{A} = \mathbf{U}$  for the PCA transformation matrix, where  $\mathbf{U}$  is the matrix of eigenvectors of  $\Sigma_{\mathbf{X}} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$ , actually diagonalizes the covariance matrix  $\Sigma_{\mathbf{Z}}$  of the transformed dataset  $\mathbf{Z}$ . Use the fact that the inverse  $\mathbf{U}^{-1} = \mathbf{U}^\top$ .

## Problem 2 (PCA for Image Analysis):

In this assignment, we apply principal component analysis (PCA) to image analysis, for the particular application of extracting eigenfaces (set of eigenvectors used extracted from images containing human faces). We begin with first setting up the environment on your local machine, then we go step by step through the procedure to extract eigenfaces.

### Setup:

- Make sure you have python3 and matplotlib installed on your system (use `pip3 install...`)
- Download the face images dataset, as well as the provided IPython notebook template from the lecture's github repository

[https://github.com/dalab/lecture\\_cil\\_public/tree/master/exercises/ex2](https://github.com/dalab/lecture_cil_public/tree/master/exercises/ex2)

### Step-by-step procedure:

1. Build a matrix collecting all images as its columns
2. Normalize all images by subtracting the mean
3. Perform PCA on the covariance matrix
4. Visualize the 5 first principal components. what is the interpretation of these images?