

Notes on Implementation of Component Analysis Techniques

Dr. Stefanos Zafeiriou

January 2015

1 Computing Principal Component Analysis

Assume that we have a matrix of centered data observations

$$\mathbf{X} = [\mathbf{x}_1 - \boldsymbol{\mu}, \dots, \mathbf{x}_N - \boldsymbol{\mu}] \quad (1)$$

where $\boldsymbol{\mu}$ denotes the mean vector. \mathbf{X} has size $F \times N$, where F is the number of dimensions and N is the number of observations. Their covariance matrix is given by

$$\mathbf{S}_t = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T = \frac{1}{N} \mathbf{X}\mathbf{X}^T \quad (2)$$

In Principal Component Analysis (PCA), we aim to maximize the variance of each dimension by maximizing

$$\begin{aligned} \mathbf{W}_0 = \arg \max_{\mathbf{W}} \quad & \text{tr}(\mathbf{W}^T \mathbf{S}_t \mathbf{W}) \\ \text{subject to} \quad & \mathbf{W}^T \mathbf{W} = \mathbf{I} \end{aligned} \quad (3)$$

The solution of Eq. 3 can be derived by solving

$$\mathbf{S}_t \mathbf{W} = \mathbf{W} \boldsymbol{\Lambda} \quad (4)$$

Thus, we need to perform eigenanalysis on \mathbf{S}_t . If we want to keep d principal components, the computational cost of the above operation is $\mathcal{O}(dF^2)$. If F is large, this computation can be quite expensive.

Lemma 1

Let us assume that $\mathbf{B} = \mathbf{X}\mathbf{X}^T$ and $\mathbf{C} = \mathbf{X}^T \mathbf{X}$. It can be proven that \mathbf{B} and \mathbf{C} have the same positive eigenvalues $\boldsymbol{\Lambda}$ and, assuming that $N < F$, then the eigenvectors \mathbf{U} of \mathbf{B} and the eigenvectors \mathbf{V} of \mathbf{C} are related as $\mathbf{U} = \mathbf{X}\mathbf{V}\boldsymbol{\Lambda}^{-\frac{1}{2}}$.

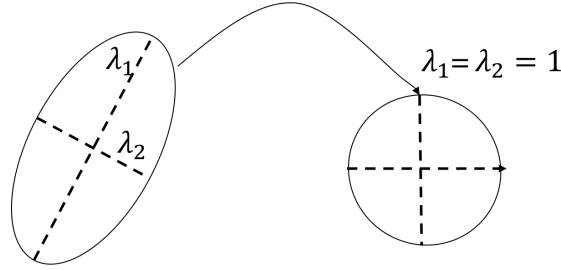


Figure 1: Example of data whitening using the PCA projection matrix $\mathbf{W} = \mathbf{U}\mathbf{\Lambda}^{-\frac{1}{2}}$.

Using Lemma 1 we can compute the eigenvectors \mathbf{U} of \mathbf{S}_t in $\mathcal{O}(N^3)$. The eigenanalysis of $\mathbf{X}^T\mathbf{X}$ is denoted by

$$\mathbf{X}^T\mathbf{X} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T \quad (5)$$

where \mathbf{V} is a $N \times (N-1)$ matrix with the eigenvectors as columns and $\mathbf{\Lambda}$ is a $(N-1) \times (N-1)$ diagonal matrix with the eigenvalues. Given that $\mathbf{V}^T\mathbf{V} = \mathbf{I}$ and $\mathbf{V}\mathbf{V}^T \neq \mathbf{I}$ we have

$$\left. \begin{aligned} \mathbf{X}^T\mathbf{X} &= \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T \\ \mathbf{U} &= \mathbf{X}\mathbf{V}\mathbf{\Lambda}^{-\frac{1}{2}} \end{aligned} \right\} \Rightarrow \mathbf{U}^T\mathbf{X}\mathbf{X}^T\mathbf{U} = \mathbf{\Lambda}^{-\frac{1}{2}}\mathbf{V}^T\mathbf{X}^T\mathbf{X}\mathbf{X}^T\mathbf{V}\mathbf{\Lambda}^{-\frac{1}{2}} =$$

$$= \mathbf{\Lambda}^{-\frac{1}{2}} \underbrace{\mathbf{V}^T\mathbf{V}}_{\mathbf{I}} \mathbf{\Lambda} \underbrace{\mathbf{V}^T\mathbf{V}}_{\mathbf{I}} \mathbf{\Lambda} \underbrace{\mathbf{V}^T\mathbf{V}}_{\mathbf{I}} \mathbf{\Lambda}^{-\frac{1}{2}} = \quad (6)$$

$$= \mathbf{\Lambda}$$

The pseudocode for computing PCA is

Algorithm 1 Principal Component Analysis

- 1: **procedure** PCA
 - 2: Compute dot product matrix: $\mathbf{X}^T\mathbf{X} = \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T (\mathbf{x}_i - \boldsymbol{\mu})$
 - 3: Eigenanalysis: $\mathbf{X}^T\mathbf{X} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$
 - 4: Compute eigenvectors: $\mathbf{U} = \mathbf{X}\mathbf{V}\mathbf{\Lambda}^{-\frac{1}{2}}$
 - 5: Keep specific number of first components: $\mathbf{U}_d = [\mathbf{u}_1, \dots, \mathbf{u}_d]$
 - 6: Compute d features: $\mathbf{Y} = \mathbf{U}_d^T\mathbf{X}$
-

Now, the covariance matrix of \mathbf{Y} is

$$\mathbf{Y}\mathbf{Y}^T = \mathbf{U}^T\mathbf{X}\mathbf{X}^T\mathbf{U} = \mathbf{\Lambda}$$

The final solution of Eq. 3 is given as the projection matrix

$$\mathbf{W} = \mathbf{U}\mathbf{\Lambda}^{-\frac{1}{2}} \quad (7)$$

which normalizes the data to have unit variance (Figure 1). This procedure is called whitening (or sphereing).

2 Computing Linear Discriminant Analysis

As explained before (Section 1), PCA finds the principal components that maximize the data variance without taking into account the class labels. In contrast to this, Linear Discriminant Analysis (LDA) computes the linear directions that maximize the separation between multiple classes. This is mathematically expressed as maximizing

$$\begin{aligned} \mathbf{W}_0 = \arg \max_{\mathbf{W}} \quad & \text{tr}(\mathbf{W}^T \mathbf{S}_b \mathbf{W}) \\ \text{subject to} \quad & \mathbf{W}^T \mathbf{S}_w \mathbf{W} = \mathbf{I} \end{aligned} \quad (8)$$

Assume that we have C number of classes, denoted by $\mathbf{c}_i = [\mathbf{x}_1, \dots, \mathbf{x}_{N_{\mathbf{c}_i}}]$, $i = 1, \dots, C$, where each \mathbf{x}_j has F dimensions and $\boldsymbol{\mu}(\mathbf{c}_i)$ is the mean vector of the class \mathbf{c}_i . Thus, the overall data matrix is $\mathbf{X} = [\mathbf{c}_1, \dots, \mathbf{c}_C]$ with size $F \times N$ ($N = \sum_{i=1}^C N_{\mathbf{c}_i}$) and $\boldsymbol{\mu}$ is the overall mean (mean of means). \mathbf{S}_w is the within-class scatter matrix

$$\mathbf{S}_w = \sum_{j=1}^C \mathbf{S}_j = \sum_{j=1}^C \sum_{\mathbf{x}_i \in \mathbf{c}_j} (\mathbf{x}_i - \boldsymbol{\mu}(\mathbf{c}_j))(\mathbf{x}_i - \boldsymbol{\mu}(\mathbf{c}_j))^T \quad (9)$$

that has $\text{rank}(\mathbf{S}_w) = \min(F, N - C)$. Moreover, \mathbf{S}_b is the between-class scatter matrix

$$\mathbf{S}_b = \sum_{j=1}^C N_{\mathbf{c}_j} (\boldsymbol{\mu}(\mathbf{c}_j) - \boldsymbol{\mu})(\boldsymbol{\mu}(\mathbf{c}_j) - \boldsymbol{\mu})^T \quad (10)$$

that has $\text{rank}(\mathbf{S}_b) = \min(F, C - 1)$. The solution of Eq. 8 is given from the generalized eigenvalue problem

$$\mathbf{S}_b \mathbf{W} = \mathbf{S}_w \mathbf{W} \boldsymbol{\Lambda} \quad (11)$$

thus \mathbf{W}_0 corresponds to the eigenvectors of $\mathbf{S}_w^{-1} \mathbf{S}_b$ that have the largest eigenvalues. In order to deal with the singularity of \mathbf{S}_w , we can do the following steps:

1. Perform PCA on our data matrix \mathbf{X} to reduce the dimensions to $N - C$ using the eigenvectors \mathbf{U}
2. Solve LDA on this reduced space and get \mathbf{Q} that has $C - 1$ columns.
3. Compute the total transform as $\mathbf{W} = \mathbf{U}\mathbf{Q}$.

Unfortunately, if you follow the above procedure is possible that important information is solved. In the following, we show how the components of LDA can be computed by applying a simultaneous diagonalization procedure.

Properties

The scatter matrices have some interesting properties. Let us denote

$$\mathbf{M} = \begin{bmatrix} \mathbf{E}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{E}_2 & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{E}_C \end{bmatrix} = \text{diag}\{\mathbf{E}_1, \mathbf{E}_2, \dots, \mathbf{E}_C\} \quad (12)$$

where

$$\mathbf{E}_i = \begin{bmatrix} \frac{1}{N_{c_i}} & \cdots & \frac{1}{N_{c_i}} \\ \vdots & \ddots & \vdots \\ \frac{1}{N_{c_i}} & \cdots & \frac{1}{N_{c_i}} \end{bmatrix}_{N_{c_i} \times N_{c_i}} \quad (13)$$

Note that \mathbf{M} is idempotent, thus $\mathbf{M}\mathbf{M} = \mathbf{M}$. Given that the data covariance matrix is $\mathbf{S}_t = \mathbf{X}\mathbf{X}^T$, the between-class scatter matrix can be written as

$$\mathbf{S}_b = \mathbf{X}\mathbf{M}\mathbf{M}\mathbf{X}^T = \mathbf{X}\mathbf{M}\mathbf{X}^T \quad (14)$$

and the within-class scatter matrix as

$$\mathbf{S}_w = \underbrace{\mathbf{X}\mathbf{X}^T}_{\mathbf{S}_t} - \underbrace{\mathbf{X}\mathbf{M}\mathbf{X}^T}_{\mathbf{S}_b} = \mathbf{X}(\mathbf{I} - \mathbf{M})\mathbf{X}^T \quad (15)$$

Thus, we have that $\mathbf{S}_t = \mathbf{S}_w + \mathbf{S}_b$. Note that since \mathbf{M} is idempotent, $\mathbf{I} - \mathbf{M}$ is also idempotent.

Given the above properties, the objective function of Eq. 8 can be expressed as

$$\begin{aligned} \mathbf{W}_0 &= \arg \max_{\mathbf{W}} \quad \text{tr}(\mathbf{W}^T \mathbf{X}\mathbf{M}\mathbf{M}\mathbf{X}^T \mathbf{W}) \\ &\text{subject to} \quad \mathbf{W}^T \mathbf{X}(\mathbf{I} - \mathbf{M})(\mathbf{I} - \mathbf{M})\mathbf{X}^T \mathbf{W} = \mathbf{I} \end{aligned} \quad (16)$$

The optimization of this problem involves a procedure called Simultaneous Diagonalization. Let's assume that the final transform matrix has the form

$$\mathbf{W} = \mathbf{U}\mathbf{Q} \quad (17)$$

We aim to find the matrix \mathbf{U} that diagonalizes $\mathbf{S}_w = \mathbf{X}(\mathbf{I} - \mathbf{M})(\mathbf{I} - \mathbf{M})\mathbf{X}^T$. This practically means that, given the constraint of Eq. 16, we want

$$\begin{aligned} \mathbf{W}^T \mathbf{X}(\mathbf{I} - \mathbf{M})(\mathbf{I} - \mathbf{M})\mathbf{X}^T \mathbf{W} &= \mathbf{I} \Rightarrow \\ \Rightarrow \mathbf{Q}^T \underbrace{\mathbf{U}^T \mathbf{X}(\mathbf{I} - \mathbf{M})(\mathbf{I} - \mathbf{M})\mathbf{X}^T \mathbf{U}}_{\mathbf{I}} \mathbf{Q} &= \mathbf{I} \end{aligned} \quad (18)$$

Consequently, using Eqs. 17 and 18, the objective function of Eq. 16 can be further expressed as

$$\begin{aligned} \mathbf{Q}_0 &= \arg \max_{\mathbf{Q}} \quad \text{tr}(\mathbf{Q}^T \mathbf{U}^T \mathbf{X}\mathbf{M}\mathbf{M}\mathbf{X}^T \mathbf{U}\mathbf{Q}) \\ &\text{subject to} \quad \mathbf{Q}^T \mathbf{Q} = \mathbf{I} \end{aligned} \quad (19)$$

where the constraint $\mathbf{W}^T \mathbf{X}(\mathbf{I} - \mathbf{M})(\mathbf{I} - \mathbf{M})\mathbf{X}^T \mathbf{W} = \mathbf{I}$ now has the form $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$.

Lemma 2

Assume the matrix $\mathbf{X}(\mathbf{I} - \mathbf{M})(\mathbf{I} - \mathbf{M})\mathbf{X}^T = \mathbf{X}_w \mathbf{X}_w^T$, where \mathbf{X}_w is the $F \times N$ matrix $\mathbf{X}_w = \mathbf{X}(\mathbf{I} - \mathbf{M})$. By performing eigenanalysis on $\mathbf{X}_w^T \mathbf{X}_w$ as $\mathbf{X}_w^T \mathbf{X}_w =$

$\mathbf{V}_w \mathbf{\Lambda} \mathbf{V}_w^T$, we get $N - C$ positive eigenvalues, thus \mathbf{V}_w is a $N \times (N - C)$ matrix.

The optimization problem of Eq. 19 can be solved in two steps

1. Find \mathbf{U} such that $\mathbf{U}^T \mathbf{X}(\mathbf{I} - \mathbf{M})(\mathbf{I} - \mathbf{M})\mathbf{X}^T \mathbf{U} = \mathbf{I}$. By applying Lemma 2, we get $\mathbf{U} = \mathbf{X}_w \mathbf{V}_w \mathbf{\Lambda}_w^{-1}$. Note that \mathbf{U} has size $F \times (N - C)$.
2. Find \mathbf{Q}_0 . By denoting

$$\tilde{\mathbf{X}}_b = \mathbf{U}^T \mathbf{X} \mathbf{M}$$

the $(N - C) \times N$ matrix of projected class means, Eq. 19 becomes

$$\begin{aligned} \mathbf{Q}_0 = \arg \max_{\mathbf{Q}} \quad & \text{tr}(\mathbf{Q}^T \tilde{\mathbf{X}}_b \tilde{\mathbf{X}}_b^T \mathbf{Q}) \\ \text{subject to} \quad & \mathbf{Q}^T \mathbf{Q} = \mathbf{I} \end{aligned} \quad (20)$$

which is equivalent to applying PCA on the matrix of projected class means. The final \mathbf{Q}_0 is a matrix with columns the d eigenvectors of $\tilde{\mathbf{X}}_b \tilde{\mathbf{X}}_b^T$ that correspond to the d largest eigenvalues ($d \leq C - 1$).

The final projection matrix is given by

$$\mathbf{W}_0 = \mathbf{Q}_0 \mathbf{U} \quad (21)$$

Based on the above, the pseudocode for computing LDA is

Algorithm 2 Linear Discriminant Analysis

- 1: **procedure** LDA
 - 2: Find eigenvectors of \mathbf{S}_w that correspond to non-zero eigenvalues (usually $N - C$), i.e. $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_{N-C}]$ by performing eigen-analysis to $(\mathbf{I} - \mathbf{M})\mathbf{X}^T \mathbf{X}(\mathbf{I} - \mathbf{M}) = \mathbf{V}_w \mathbf{\Lambda}_w \mathbf{V}_w^T$ and computing $\mathbf{U} = \mathbf{X}(\mathbf{I} - \mathbf{M})\mathbf{V}_w \mathbf{\Lambda}_w^{-1}$ (performing whitening on \mathbf{S}_w).
 - 3: Project the data as $\tilde{\mathbf{X}}_b = \mathbf{U}^T \mathbf{X} \mathbf{M}$.
 - 4: Perform PCA on $\tilde{\mathbf{X}}_b$ to find \mathbf{Q} (i.e., compute the eigenanalysis of $\tilde{\mathbf{X}}_b \tilde{\mathbf{X}}_b^T = \mathbf{Q} \mathbf{\Lambda}_b \mathbf{Q}^T$).
 - 5: The total transform is $\mathbf{W} = \mathbf{U} \mathbf{Q}$
-

Locality preserving projections can be computed in a similar fashion. The first step is to perform whitening of $\mathbf{X} \mathbf{D} \mathbf{X}^T$. We do so by applying Lemma 1 and performing eigenanalysis of $\mathbf{D}^{1/2} \mathbf{X}^T \mathbf{X} \mathbf{D}^{1/2} = \mathbf{V}_p \mathbf{\Lambda}_p \mathbf{V}_p^T$. Then, the whitening transform is given by $\mathbf{U}^T = \mathbf{X} \mathbf{D}^{1/2} \mathbf{\Lambda}_p^{-1/2}$. The next step is to project the data as $\tilde{\mathbf{X}}_p = \mathbf{U}^T \mathbf{X}$ and find the eigenvectors \mathbf{Q} of $\tilde{\mathbf{X}}_p (\mathbf{D} - \mathbf{S}) \tilde{\mathbf{X}}_p^T$ that correspond to the lowest (but non-zero eigenvalues). Then the total transform would be $\mathbf{W} = \mathbf{U} \mathbf{Q}$.

¹Where \mathbf{S} is the connectivity matrix defined in the slides