

PCA Derivation

Fall 2018

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

1 Problem Statement

Define $\mu, \mathbf{A}, \theta_1, \dots, \theta_n$ to be the solution of

$$\min_{\substack{\mu \in \mathbb{R}^d \\ \mathbf{A} \in \mathcal{A}_k \\ \theta_i \in \mathbb{R}^k}} \sum_{i=1}^n \|\mathbf{x}_i - \mu - \mathbf{A}\theta_i\|^2$$

where \mathcal{A}_k is the set of $d \times k$ matrices with orthonormal columns. PCA gives the least squares rank- k linear approximation to the data set. The solution to above optimization problem is given in terms of the spectral (or eigenvalue) decomposition of the sample covariance matrix

$$S = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

In particular, write

$$S = U \Lambda U^T$$

where

$$U = [\mathbf{u}_1 \cdots \mathbf{u}_d] \text{ and } \lambda = \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_d \end{bmatrix}$$

with $U^T U = U U^T = I$, $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \geq 0$. Recall that $S \mathbf{u}_j = \lambda_j \mathbf{u}_j \forall j$. We will show that one solution to above optimization problem is

$$\begin{aligned} \mu &= \bar{\mathbf{x}} \\ \mathbf{A} &= [\mathbf{u}_1 \cdots \mathbf{u}_k] \\ \theta_i &= \mathbf{A}^T (\mathbf{x}_i - \bar{\mathbf{x}}). \end{aligned}$$

We will also characterize the set of all solutions. Some terminology:

- principal component transform:

$$\mathbf{x} \mapsto \mathbf{A}^T (\mathbf{x} - \bar{\mathbf{x}}) \in \mathbb{R}^k$$

- j^{th} principal component:

$$\theta^{(j)} = \mathbf{u}_j^T (\mathbf{x} - \bar{\mathbf{x}}) \in \mathbb{R}$$

- j^{th} principal eigenvector:

$$\mathbf{u}_j \in \mathbb{R}^d$$

2 Overview

We want to minimize

$$\sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu} - A\boldsymbol{\theta}_i\|^2 \quad (1)$$

with respect to $\boldsymbol{\mu} \in \mathbb{R}^d$, $A \in \mathcal{A}_k$, $\boldsymbol{\theta}_i \in \mathbb{R}^k$

Step 1: Eliminate $\boldsymbol{\theta}_i$

Suppose A , $\boldsymbol{\mu}$ are fixed. We can optimize each term with respect to $\boldsymbol{\theta}_i$ individually, yielding

$$\begin{aligned} \boldsymbol{\theta}_i &= (A^T A)^{-1} A^T (\mathbf{x}_i - \boldsymbol{\mu}) \\ &= A^T (\mathbf{x}_i - \boldsymbol{\mu}) \end{aligned} \quad (2)$$

Step 2: Eliminate $\boldsymbol{\mu}$

Holding A fixed, we wish to minimize

$$\begin{aligned} &\sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu} - AA^T(\mathbf{x}_i - \boldsymbol{\mu})\|^2 \\ &= \sum_{i=1}^n \|(I - AA^T)(\mathbf{x}_i - \boldsymbol{\mu})\|^2 \\ &= \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^T (I - AA^T)^T (I - AA^T) (\mathbf{x}_i - \boldsymbol{\mu}) \end{aligned} \quad (3)$$

Note that $I - AA^T$ is a projection matrix onto the orthogonal complement of $\langle A \rangle$. Therefore

$$\begin{aligned} (I - AA^T)^T (I - AA^T) &= (I - AA^T)(I - AA^T) \\ &= I - AA^T \end{aligned} \quad (4)$$

Note that $I - AA^T$, being a projection matrix, is PSD, and therefore

$$\sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^T (I - AA^T) (\mathbf{x}_i - \boldsymbol{\mu}) \quad (5)$$

is a convex function of $\boldsymbol{\mu}$. Now

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\mu}} &= 2 \sum_{i=1}^n (I - AA^T)(\mathbf{x}_i - \boldsymbol{\mu}) \\ &= 2(I - AA^T) \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}) \\ &= 2n(I - AA^T)(\bar{\mathbf{x}} - \boldsymbol{\mu}) = \mathbf{0} \end{aligned} \quad (6)$$

is solved by $\boldsymbol{\mu} = \bar{\mathbf{x}}$. More generally, it suffices for $\bar{\mathbf{x}} - \boldsymbol{\mu}$ to belong to the nullspace of $I - AA^T$, which is $\langle A \rangle$. Thus any $\boldsymbol{\mu} \in \bar{\mathbf{x}} + \langle A \rangle$ is a possible solution.

Step 3: Optimize A

It remains to solve

$$\min_{A \in \mathcal{A}_k} \sum \|\mathbf{x}_i - \bar{\mathbf{x}} - AA^T(\mathbf{x}_i - \bar{\mathbf{x}})\|^2 \quad (7)$$

Assume $\bar{\mathbf{x}} = \mathbf{0}$ which can be ensured by subtracting $\bar{\mathbf{x}}$ from each \mathbf{x}_i . Also note that AA^T is a rank k projection matrix. Let \mathcal{P}_k denote the set of $d \times d$ rank- k projection matrices. Then it remains to solve

$$\min_{P \in \mathcal{P}_k} \sum \|\mathbf{x}_i - P\mathbf{x}_i\|^2 \quad (8)$$

Introduce the data matrix

$$X = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_n] \quad (d \times n) \quad (9)$$

and for an arbitrary matrix $C = [c_{ij}]$, define the Frobenius norm

$$\|C\|_F := \sqrt{\sum_i \sum_j c_{ij}^2} \quad (10)$$

Then we can restate the problem as

$$\min_{P \in \mathcal{P}_k} \|X - PX\|_F^2. \quad (11)$$

This is the core optimization problem at the heart of PCA. Although it is nonconvex, it is a very special nonconvex problem and has a closed form solution given by $P = AA^T$ where A consists of the first k principal eigenvectors. The derivation of this result is quite interesting in its own right, and connects with other important topics in matrix algebra. These details are given in the next section.

3 The Fundamental PCA Problem

Above we saw that PCA reduces to the following optimization problem:

$$\min_{P \in \mathcal{P}_k} \|\mathbf{X} - \mathbf{P}\mathbf{X}\|_F \quad (\text{PCA})$$

where \mathbf{X} is the $d \times n$ data matrix (column mean = 0), $\|\cdot\|_F$ is the Frobenius norm, and \mathcal{P}_k is the set of rank k , $d \times d$ projection matrices. In this setting, the covariance matrix of the data is

$$S = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{n} \mathbf{X} \mathbf{X}^T$$

(verify by comparing entries).

We will derive the solution to PCA in two different ways, by connecting (PCA) to two other problems in matrix algebra.

4 Connection to the SVD

Every matrix \mathbf{X} has a singular value decomposition (SVD)

$$\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$$

where the columns of \mathbf{U} are the left singular vectors, the columns of \mathbf{V} are the right singular vectors, $\sigma_1 \geq \sigma_2 \geq \cdots \geq 0$ are singular values and

- $\mathbf{U} \mathbf{U}^T = \mathbf{U}^T \mathbf{U} = \mathbf{I}_{d \times d}$
- $\mathbf{V} \mathbf{V}^T = \mathbf{V}^T \mathbf{V} = \mathbf{I}_{n \times n}$

$$\bullet \Sigma = \begin{cases} \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_d \end{bmatrix} & \text{if } n \geq d \\ \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_d \end{bmatrix} & \text{if } n < d \end{cases}$$

The SVD arises in the following theorem due to Eckart and Young.

Theorem 1. *Let \mathbf{X} have rank $r = k$. The solution to*

$$\min_{\mathbf{Z} \in \mathbb{R}^{d \times n}, \text{rank}(\mathbf{Z})=k} \|\mathbf{X} - \mathbf{Z}\|_F^2 \quad (\text{SVD})$$

is $\mathbf{Z}_k = \mathbf{U}\Sigma_k\mathbf{V}^T$, where Σ_k is Σ with $\sigma_{k+1}, \sigma_{k+2}, \dots$ set to zero.

To see the connection to (PCA), write

$$\Sigma_k = \mathbf{I}_{k,d}\Sigma,$$

where

$$\mathbf{I}_{k,d} = \begin{bmatrix} 1 & & & & & \\ & 1 & & & & \\ & & \ddots & & & \\ & & & 1 & & \\ & & & & 0 & \\ & & & & & 0 \\ & & & & & & \ddots & \\ & & & & & & & 0 \end{bmatrix}$$

with k 1's, in a $(d \times d)$ matrix. Then

$$\begin{aligned} \mathbf{Z}_k &= \mathbf{U}\Sigma_k\mathbf{V}^T \\ &= \mathbf{U}(\mathbf{I}_{k,d}\Sigma)\mathbf{V}^T \\ &= \mathbf{U}(\mathbf{I}_{k,d}\mathbf{U}^T\mathbf{X}\mathbf{V})\mathbf{V}^T \\ &= \mathbf{U}_k\mathbf{U}_k^T\mathbf{X} \end{aligned}$$

where \mathbf{U}_k contains the first k left singular vectors. Clearly $\mathbf{U}_k\mathbf{U}_k^T \in \mathbf{P}_k$. Therefore $\mathbf{P} = \mathbf{U}_k\mathbf{U}_k^T$ gives a solution to (PCA). It remains to show that the left singular vectors are the eigenvectors of $\mathbf{X}\mathbf{X}^T$. To see this, observe that the eigenvalue decomposition of $\mathbf{X}\mathbf{X}^T$ is

$$\begin{aligned} \mathbf{X}\mathbf{X}^T &= \mathbf{U}\Sigma\mathbf{V}^T\mathbf{V}\Sigma^T\mathbf{U}^T \\ &= \mathbf{U}\Sigma\Sigma^T\mathbf{U}^T \\ &= \mathbf{U}\Lambda\mathbf{U}^T \end{aligned}$$

where

$$\mathbf{\Lambda} = \begin{cases} \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_d^2 \end{bmatrix} & \text{if } n \geq d \\ \begin{bmatrix} \sigma_1^2 & & & & \\ & \ddots & & & \\ & & \sigma_n^2 & & \\ & & & 0 & \\ & & & & \ddots \\ & & & & & 0 \end{bmatrix} & \text{if } n < d \end{cases}$$

Therefore, the SVD of \mathbf{X} gives the (PCA) solution:

- principal eigenvectors = left singular vectors of \mathbf{X}
- $\lambda_i = \frac{1}{n} (i^{th} \text{ singular value of } \mathbf{X})^2$

5 Generalized Rayleigh Quotient

We now present a second solution of (PCA).

The trace of a square matrix is the sum of the diagonal entries. It satisfies the following properties:

- linearity: $tr(\mathbf{C} + \mathbf{D}) = tr(\mathbf{C}) + tr(\mathbf{D})$ for any two square matrices \mathbf{C} and \mathbf{D}
- invariance to cyclic permutations: $tr(\mathbf{CD}) = tr(\mathbf{DC})$
- the trace of a symmetric matrix is the sum of its eigenvalues
- for any matrix \mathbf{C} , $\|\mathbf{C}\|_F^2 = tr(\mathbf{C}^T \mathbf{C})$

These properties are easily verified. Now observe

$$\begin{aligned} \|\mathbf{X} - \mathbf{PX}\|_F^2 &= tr((\mathbf{X} - \mathbf{PX})^T (\mathbf{X} - \mathbf{PX})) \\ &= tr(\mathbf{X}^T \mathbf{X}) - tr(\mathbf{X}^T \mathbf{PX}) - tr(\mathbf{X}^T \mathbf{P}^T \mathbf{X}) + tr(\mathbf{X}^T \mathbf{P}^T \mathbf{PX}) \\ &= tr(\mathbf{X}^T \mathbf{X}) - tr(\mathbf{X}^T \mathbf{PX}) \end{aligned}$$

where we used $\mathbf{P} = \mathbf{P}^T$ and $\mathbf{P}^2 = \mathbf{P}$. Writing $\mathbf{P} = \mathbf{AA}^T$ where $\mathbf{A} \in \mathcal{A}_k$, we need to maximize

$$\begin{aligned} tr(\mathbf{X}^T \mathbf{PX}) &= tr(\mathbf{X}^T \mathbf{AA}^T \mathbf{X}) \\ &= tr(\mathbf{A}^T \mathbf{XX}^T \mathbf{A}) \end{aligned}$$

The derivation of PCA is concluded by the following result, which I will refer to as the Generalized Rayleigh Quotient theorem (this terminology is not standard) because, in the special case $k = 1$, it relates to the Rayleigh quotient

Theorem 2. Let \mathbf{C} be a PSD matrix with eigenvalue decomposition $\mathbf{C} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$, where $\mathbf{U} = [\mathbf{u}_1 \cdots \mathbf{u}_d]$. Then a solution of

$$\max_{\mathbf{A} \in \mathcal{A}_K} tr(\mathbf{A}^T \mathbf{CA}) \tag{GRQ}$$

is $\mathbf{A} = [\mathbf{u}_1 \cdots \mathbf{u}_k]$.

6 Proof of Generalized Rayleigh Quotient Theorem

Introduce the change of variable

$$\mathbf{w}_i = \mathbf{U}^T \mathbf{a}_i$$

We know that $\mathbf{w}_1, \dots, \mathbf{w}_k$ are orthonormal because

$$\mathbf{w}_i^T \mathbf{w}_j = \mathbf{a}_i^T \mathbf{U} \mathbf{U}^T \mathbf{a}_j = \mathbf{a}_i^T \mathbf{a}_j$$

Then we need to maximize

$$\begin{aligned} \text{tr}(\mathbf{A}^T \mathbf{C} \mathbf{A}) &= \sum_{i=1}^k \mathbf{a}_i^T \mathbf{C} \mathbf{a}_i \\ &= \sum_{i=1}^k \mathbf{w}_i^T \mathbf{\Lambda} \mathbf{w}_i \end{aligned}$$

subject to $\mathbf{W} = [\mathbf{w}_1 \cdots \mathbf{w}_k] \in \mathcal{A}_k$. Now

$$\begin{aligned} \sum_{i=1}^k \mathbf{w}_i^T \mathbf{\Lambda} \mathbf{w}_i &= \sum_{i=1}^k \sum_{j=1}^d (w_i^{(j)})^2 \lambda_j \\ &= \sum_{j=1}^d \left[\sum_{i=1}^k (w_i^{(j)})^2 \right] \lambda_j \\ &= \sum_{j=1}^d h_j \lambda_j \end{aligned}$$

where $h_j = \sum_{i=1}^k (w_i^{(j)})^2$

Lemma 1. $0 \leq h_j \leq 1 \ \forall j$ and $\sum_{j=1}^d h_j = k$.

Proof. The second part is easy:

$$\begin{aligned} \sum_{j=1}^d h_j &= \sum_{j=1}^d \left(\sum_{i=1}^k (w_i^{(j)})^2 \right) \\ &= \sum_{i=1}^k \left(\sum_{j=1}^d (w_i^{(j)})^2 \right) \\ &= \sum_{i=1}^k (1) \\ &= k \end{aligned}$$

$h_j \geq 0$ is also obvious. To show $h_j \leq 1$, let $\mathbf{w}_{k+1}, \dots, \mathbf{w}_d$ extend $\mathbf{w}_1, \dots, \mathbf{w}_k$ to an orthonormal basis. Consider the $d \times d$ matrix

$$\mathbf{M} = [\mathbf{w}_1 \cdots \mathbf{w}_d] \tag{12}$$

We know $\mathbf{M}^T \mathbf{M} = \mathbf{I}$ by orthonormality. Therefore \mathbf{M}^T is a right inverse of \mathbf{M} , and so must also be a right inverse (a property of square matrices), meaning $\mathbf{M} \mathbf{M}^T = \mathbf{I}$. This implies

$$h_j = \sum_{i=1}^k (w_i^{(j)})^2 \leq \sum_{i=1}^d (w_i^{(j)})^2 = 1$$

□

We need to maximize

$$\sum_{j=1}^d h_j \lambda_j$$

with respect to the constraints imposed by the lemma.

Since $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$, this is accomplished by

$$h_j = \begin{cases} 1 & \text{if } 1 \leq j \leq k \\ 0 & \text{otherwise} \end{cases}$$

which in turn is achieved by

$$\mathbf{W} = \begin{bmatrix} I_{k \times k} \\ \hline 0 \end{bmatrix}$$

Therefore $\mathbf{A} = \mathbf{U}\mathbf{W} = [\mathbf{u}_1 \dots \mathbf{u}_k]$. Note that the optimal \mathbf{A} is not unique. Indeed if

$$\mathbf{W} = \begin{bmatrix} \text{any set of length } k \\ \text{orthonormal vectors} \\ \hline 0 \end{bmatrix}$$

then $\mathbf{A} = \mathbf{U}\mathbf{W}$ also achieves the maximum.

An interesting question is when is $\langle \mathbf{A} \rangle$ unique. This is left as an exercise.

The GRQ theorem can be used to derive PCA from the sequential maximum variance approach discussed in class. Furthermore, instead of a sequential definition of maximum variance, we can instead ask what $\mathbf{A} \in \mathcal{A}_k$ maximizes the **total variance**

$$\sum_{i=1}^k \mathbf{a}_i^T \mathbf{S} \mathbf{a}_i = \text{tr}(\mathbf{A}^T \mathbf{S} \mathbf{A}).$$

The GRQ theorem again tells us that the principal eigenvectors of \mathbf{S} provide the solution. We will employ the GRQ theorem again later in the course when we discuss spectral clustering.