

Derivation of Principal Components Analysis

Given a set of data points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, we want to find the linear subspace (plus an affine offset) that is the best fit in the least-squares sense. Mathematically, we want to solve

$$\underset{\boldsymbol{\mu}, \mathbf{A}, \{\boldsymbol{\theta}_i\}}{\text{minimize}} \quad \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{A}\boldsymbol{\theta}_i\|_2^2, \quad \text{subject to } \mathbf{A}^T \mathbf{A} = \mathbf{I},$$

where $\boldsymbol{\mu} \in \mathbb{R}^d$, $\boldsymbol{\theta}_i \in \mathbb{R}^k$, and \mathbf{A} is a $n \times k$ matrix; the constraint $\mathbf{A}^T \mathbf{A} = \mathbf{I}$ means that we will consider \mathbf{A} with orthonormal columns.

Minimizing the expression above over the $\{\boldsymbol{\theta}_i\}$ and $\boldsymbol{\mu}$ is straightforward. For the $\{\boldsymbol{\theta}_i\}$, suppose that \mathbf{A} and $\boldsymbol{\mu}$ are fixed. Then we have a series of n decoupled least-squares problems: for $i = 1, \dots, n$, we solve

$$\underset{\boldsymbol{\theta}_i}{\text{minimize}} \quad \|\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{A}\boldsymbol{\theta}_i\|_2^2$$

This is a standard unconstrained least-squares problem that has solution

$$\hat{\boldsymbol{\theta}}_i = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{x}_i - \boldsymbol{\mu}) = \mathbf{A}^T (\mathbf{x}_i - \boldsymbol{\mu}),$$

where the second equality follows from $\mathbf{A}^T \mathbf{A} = \mathbf{I}$. With \mathbf{A} still fixed, we solve for $\boldsymbol{\mu}$ by plugging in our expression for the $\boldsymbol{\theta}_i$:

$$\begin{aligned} \underset{\boldsymbol{\mu}}{\text{minimize}} \quad & \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{A}\mathbf{A}^T(\mathbf{x}_i - \boldsymbol{\mu})\|_2^2, \\ & = \sum_{i=1}^n \|(\mathbf{I} - \mathbf{A}\mathbf{A}^T)(\mathbf{x}_i - \boldsymbol{\mu})\|_2^2, \\ & = \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^T (\mathbf{I} - \mathbf{A}\mathbf{A}^T) (\mathbf{x}_i - \boldsymbol{\mu}), \end{aligned}$$

where the last step comes from expanding out the norm squared as an inner product, and using the fact that $(\mathbf{I} - \mathbf{A}\mathbf{A}^T)$ is a projector; it is symmetric, and $(\mathbf{I} - \mathbf{A}\mathbf{A}^T)^2 = (\mathbf{I} - \mathbf{A}\mathbf{A}^T)$. Taking the gradient of the expression above and setting it to zero means that $\hat{\boldsymbol{\mu}}$ will obey

$$\begin{aligned} \mathbf{0} &= -2 \sum_{i=1}^n (\mathbf{I} - \mathbf{A}\mathbf{A}^T)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}) \\ &= -2(\mathbf{I} - \mathbf{A}\mathbf{A}^T) \left(\sum_{i=1}^n \mathbf{x}_i - n\hat{\boldsymbol{\mu}} \right). \end{aligned}$$

This can be satisfied by taking

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i.$$

Note that this is not the only choice for $\boldsymbol{\mu}$ — any choice that puts $\sum_i \mathbf{x}_i - n\boldsymbol{\mu}$ into the column space of \mathbf{A} will work. But the choice above is intuitive, so we will go with it.

With $\{\hat{\boldsymbol{\theta}}_i\}$ and $\hat{\boldsymbol{\mu}}$ solved for, we now optimize over \mathbf{A} . We want to solve

$$\underset{\mathbf{A} \in \mathbb{R}^{n \times k}}{\text{minimize}} \quad \sum_{i=1}^n \|\mathbf{x}_i - \hat{\boldsymbol{\mu}} - \mathbf{A}\mathbf{A}^T(\mathbf{x}_i - \hat{\boldsymbol{\mu}})\|_2^2 \quad \text{subject to} \quad \mathbf{A}^T \mathbf{A} = \mathbf{I}.$$

We will assume, without loss of generality, that $\hat{\boldsymbol{\mu}} = \mathbf{0}$, as we could simply use the variable substitution $\tilde{\mathbf{x}}_i = \mathbf{x}_i - \hat{\boldsymbol{\mu}}$ above. The program becomes

$$\underset{\mathbf{A} \in \mathbb{R}^{n \times k}}{\text{minimize}} \quad \sum_{i=1}^n \|(\mathbf{I} - \mathbf{A}\mathbf{A}^T)\mathbf{x}_i\|_2^2 \quad \text{subject to} \quad \mathbf{A}^T \mathbf{A} = \mathbf{I}.$$

Expanding the functional, and again using the fact that $(\mathbf{I} - \mathbf{A}\mathbf{A}^T)$ is a projector,

$$\begin{aligned} \sum_{i=1}^n \|(\mathbf{I} - \mathbf{A}\mathbf{A}^T)\mathbf{x}_i\|_2^2 &= \sum_{i=1}^n \mathbf{x}_i^T (\mathbf{I} - \mathbf{A}\mathbf{A}^T) \mathbf{x}_i \\ &= \sum_{i=1}^n \mathbf{x}_i^T \mathbf{x}_i - \mathbf{x}_i^T \mathbf{A}\mathbf{A}^T \mathbf{x}_i. \end{aligned}$$

The first term does not depend on \mathbf{A} , and the second term is always negative, so our problem is equivalent to

$$\underset{\mathbf{A} \in \mathbb{R}^{n \times k}}{\text{maximize}} \quad \sum_{i=1}^n \mathbf{x}_i^T \mathbf{A}\mathbf{A}^T \mathbf{x}_i \quad \text{subject to} \quad \mathbf{A}^T \mathbf{A} = \mathbf{I}.$$

For any vector \mathbf{v} , it is easy to see that $\|\mathbf{v}\|_2^2 = \text{trace}(\mathbf{v}\mathbf{v}^T)$. Thus, the objective function above can also be written as

$$\begin{aligned} \sum_{i=1}^n \mathbf{x}_i^T \mathbf{A}\mathbf{A}^T \mathbf{x}_i &= \sum_{i=1}^n \|\mathbf{A}^T \mathbf{x}_i\|_2^2 \\ &= \sum_{i=1}^n \text{trace}(\mathbf{A}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{A}) \\ &= \text{trace} \left(\mathbf{A}^T \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{A} \right) \\ &= \text{trace}(\mathbf{A}^T \mathbf{S} \mathbf{A}), \end{aligned}$$

where $\mathbf{S} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$ is a scaled version of the sample covariance matrix.

By construction, \mathbf{S} is symmetric positive semi-definite, so it has eigenvalue decomposition

$$\mathbf{S} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T,$$

where \mathbf{U} is a $d \times d$ orthonormal matrix, $\mathbf{U}^T \mathbf{U} = \mathbf{U} \mathbf{U}^T = \mathbf{I}$, and $\mathbf{\Lambda} = \text{diag}(\{\lambda_i\})$, with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$. Then

$$\text{trace}(\mathbf{A}^T \mathbf{S} \mathbf{A}) = \text{trace}(\mathbf{A}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \mathbf{A}) = \text{trace}(\mathbf{W}^T \mathbf{\Lambda} \mathbf{W}),$$

where $\mathbf{W} = \mathbf{U}^T \mathbf{A}$. Notice that \mathbf{W} also has orthonormal columns, as $\mathbf{W}^T \mathbf{W} = \mathbf{A}^T \mathbf{U} \mathbf{U}^T \mathbf{A} = \mathbf{A}^T \mathbf{A} = \mathbf{I}$. So we can solve the program

$$\underset{\mathbf{W} \in \mathbb{R}^{n \times k}}{\text{maximize}} \text{trace}(\mathbf{W}^T \mathbf{\Lambda} \mathbf{W}) \quad \text{subject to} \quad \mathbf{W}^T \mathbf{W} = \mathbf{I},$$

and then take $\hat{\mathbf{A}} = \mathbf{U}^T \hat{\mathbf{W}}$.

We can show that the last maximization program above is equivalent to a simple linear program that we can solve by inspection. Let $\mathbf{w}_1, \dots, \mathbf{w}_k$ be the columns of \mathbf{W} . Then

$$\begin{aligned} \text{trace}(\mathbf{W}^T \mathbf{\Lambda} \mathbf{W}) &= \sum_{i=1}^k \mathbf{w}_i^T \mathbf{\Lambda} \mathbf{w}_i \\ &= \sum_{i=1}^k \sum_{j=1}^d w_i(j)^2 \lambda_j \\ &= \sum_{j=1}^d h_j \lambda_j, \quad h_j = \sum_{i=1}^k w_i(j)^2 = \sum_{i=1}^k W(i, j)^2. \end{aligned}$$

The h_j , $j = 1, \dots, d$ above are the sums of the squares of the *rows* of \mathbf{W} . It is clear that $h_j \geq 0$. It is also true that

$$\sum_{j=1}^d h_j = k,$$

as the fact that the norm of each columns of \mathbf{W} is 1 means that

$$\sum_{j=1}^d \sum_{i=1}^k W(i, j)^2 = \sum_{i=1}^k \left(\sum_{j=1}^d W(i, j)^2 \right) = \sum_{i=1}^k 1 = k.$$

Finally, it is also true that $h_j \leq 1$. Here is why: since the columns of \mathbf{W} are orthonormal, they can be considered as part of an orthonormal basis for all of \mathbb{R}^d . That is, there is a (and actually there are many) $d \times (d - k)$ matrix \mathbf{W}_0 such that the columns of

$$\mathbf{W}' = [\mathbf{W} \quad \mathbf{W}_0]$$

form an orthonormal basis for \mathbb{R}^d . Since \mathbf{W}' is square, $\mathbf{W}'\mathbf{W}'^T = \mathbf{I}$, meaning the sum of the squares of each row are equal to 1. Thus

$$h_j = \sum_{i=1}^k W(i, j)^2 \leq \sum_{i=1}^d W'(i, j)^2 = 1.$$

With these constraints on the h_j , let's see how large we can make the quantity of interest:

$$\underset{\mathbf{h} \in \mathbb{R}^d}{\text{maximize}} \quad \sum_{j=1}^d h_j \lambda_j \quad \text{subject to} \quad \sum_{j=1}^d h_j = k, \quad 0 \leq h_j \leq 1.$$

This is a linear program, but we can intuit the answer. Since all of the λ_j are positive, we want to have their weights (i.e., the h_j) as large as possible for the largest entries. Since the weights are constrained to be less than 1, and their sum is k , this simply means we assign a weight of 1 to the k largest terms, and 0 to the others:

$$\widehat{h}_j = \begin{cases} 1, & j = 1, \dots, k, \\ 0, & \text{otherwise.} \end{cases}$$

This means that the sum of the squares of the entries in the rows of the corresponding $\widehat{\mathbf{W}}$ are 1 for the first k , and zero below — there

are many matrices with orthonormal columns which fit the bill, but a specific one which does is

$$\widehat{\mathbf{W}} = \begin{bmatrix} \mathbf{I}_{k \times k} \\ \mathbf{0}_{(d-k) \times k} \end{bmatrix}. \quad (1)$$

Taking $\widehat{\mathbf{A}} = \mathbf{U}^T \widehat{\mathbf{W}}$, this results in

$$\widehat{\mathbf{A}} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_k],$$

where the \mathbf{u}_i above are the first k columns of \mathbf{U} .

PCA Theorem

$$\underset{\mu, \mathbf{A}, \{\theta_i\}}{\text{minimize}} \sum_{i=1}^n \|\mathbf{x}_i - \mu - \mathbf{A}\theta_i\|_2^2, \quad \text{subject to } \mathbf{A}^T \mathbf{A} = \mathbf{I},$$

has solution

$$\widehat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad \widehat{\mathbf{A}} = [\mathbf{u}_1 \ \cdots \ \mathbf{u}_k], \quad \widehat{\theta}_i = \widehat{\mathbf{A}}^T (\mathbf{x}_i - \widehat{\mu}),$$

where $\mathbf{u}_1, \dots, \mathbf{u}_k$ are the eigenvectors corresponding to the k largest eigenvalues of

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

Note that our analysis above shows that the choice of \mathbf{A} is not unique — we are really choosing the subspace spanned by the columns of \mathbf{A} , and do not care which orthobasis we use to span it. In the end,

taking $\widehat{\mathbf{A}}' = \widehat{\mathbf{A}}\mathbf{Q}$, for any $k \times k$ orthonormal matrix \mathbf{Q} would also work, as

$$\widehat{\mathbf{A}}'\widehat{\mathbf{A}}'^{\text{T}} = \widehat{\mathbf{A}}\mathbf{Q}\mathbf{Q}^{\text{T}}\widehat{\mathbf{A}}^{\text{T}} = \widehat{\mathbf{A}}\widehat{\mathbf{A}}^{\text{T}}.$$

In our choice for $\widehat{\mathbf{W}}$ in (1) above, we would take

$$\widehat{\mathbf{W}} = \begin{bmatrix} \mathbf{Q} \\ \mathbf{0}_{(d-k) \times k} \end{bmatrix},$$

which also meets the constraints dictated by the \widehat{h}_j — the sum of the squares of the entries in the rows is 1 for the first k , zero for the last $d - k$, and the columns are orthonormal.