

1 Singular Value Decomposition and Principal Component Analysis

In these lectures we discuss the SVD and the PCA, two of the most widely used tools in machine learning. Principal Component Analysis (PCA) is a linear dimensionality reduction method dating back to Pearson (1901) and it is one of the most useful techniques in exploratory data analysis. It is also known under different names such as the Karhunen-Love Transform, the Hotelling transform, and Proper Orthogonal Decomposition (POD). PCA can be applied to a data set comprising of n vectors $x_1, \dots, x_n \in \mathbb{R}^d$ and in turn returns a new basis for \mathbb{R}^d whose elements are terms the principal components. It is important that the method is completely data-dependent, that is, the new basis is only a function of the data. The PCA builds on the SVD (or the spectral theorem), we therefore start with the SVD.

1.1 Singular Value Decomposition (SVD)

Consider a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ or $\mathbb{C}^{m \times n}$ and let us assume that $m \geq n$. Then the *singular value decomposition* (SVD) of \mathbf{A} is given by [1]

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{W}^*,$$

where \mathbf{U} is $m \times m$, \mathbf{D} is $m \times n$, \mathbf{W} is $n \times n$, \mathbf{U} and \mathbf{W} are unitary (i.e., $\mathbf{U}^* \mathbf{U} = \mathbf{U} \mathbf{U}^* = \mathbf{I}_m$, $\mathbf{W} \mathbf{W}^* = \mathbf{W}^* \mathbf{W} = \mathbf{I}_n$), and \mathbf{D} is a diagonal (rectangular) matrix

$$\mathbf{D} = \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \sigma_n \\ 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

with $D_{ii} = \sigma_i > 0$. Here, σ_i are called the *singular values* of \mathbf{A} , the columns of \mathbf{U} are the corresponding *left singular vectors*, and the columns of \mathbf{W} are the corresponding *right singular vectors*.

Let $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$, $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_n]$ and let r be the rank of \mathbf{A} . Then we can write

$$\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{w}_i^*,$$

with $r \leq n$ (and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$). (So \mathbf{A} is a sum of weighted rank-one matrices.) The SVD exists for any finite-dimensional matrix.

Remarks:

- The \mathbf{u}_i are eigenvectors of $\mathbf{A}\mathbf{A}^*$ and the \mathbf{w}_i are eigenvectors of $\mathbf{A}^*\mathbf{A}$.
- $\mathbf{A}\mathbf{A}^*$ and $\mathbf{A}^*\mathbf{A}$ are positive semidefinite so their eigenvalues are nonnegative.
- If λ_i are the eigenvalues of $\mathbf{A}^*\mathbf{A}$, then $\sigma_i^2 = \lambda_i$ if $\lambda_i > 0$. (Here we're saying that singular values must be positive, but this is more of a matter of taste.)
- If \mathbf{A} is square and Hermitian, then the SVD and the eigenvalue decomposition are the same.
- We could alternatively define the SVD with \mathbf{U} as an $m \times n$ matrix, \mathbf{D} as an $n \times n$ matrix, and \mathbf{W} as an $n \times n$ matrix. In this case, $\mathbf{U}^*\mathbf{U} = \mathbf{I}_n$, and $\mathbf{W}^*\mathbf{W} = \mathbf{W}\mathbf{W}^* = \mathbf{I}_n$.

Some intuition for SVD: SVD rotates the matrix \mathbf{A} by \mathbf{U} and \mathbf{W}^* so that \mathbf{A} becomes a diagonal matrix.

2 Principal Component Analysis (PCA)

2.1 Motivation

Given $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, we want to project the \mathbf{x}_i onto \mathbb{R}^k , $k < d$. So, how do we choose k and the orientation of the subspace? We consider two ideas:

1. Find the k -dimensional subspace for which the projections of $\mathbf{x}_1, \dots, \mathbf{x}_n$ best approximate the original points $\mathbf{x}_1, \dots, \mathbf{x}_n$. (We define “best approximation” in the sense of the 2-norm.)
2. We also want to conserve what makes the data points different from each other. Hence, find the k -dimensional projection of $\mathbf{x}_1, \dots, \mathbf{x}_n$ that preserves most of the variance of the \mathbf{x}_i .

Both of the two ideas above are solved by *principal component analysis* (PCA).

2.2 Optimization problem formulation

We denote the *sample mean* by

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

and *sample covariance* matrix by

$$\Sigma_n := \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}_n) (\mathbf{x}_i - \boldsymbol{\mu}_n)^*.$$

Let us focus on the first idea in Section 2.1. We want to approximate each \mathbf{x}_i by an affine low-dimensional subspace such that for each \mathbf{x}_i we have

$$\mathbf{x}_i \approx \boldsymbol{\mu} + \sum_{j=1}^k (\boldsymbol{\alpha}_i)_j \mathbf{v}_j,$$

where $\mathbf{V} := [\mathbf{v}_1, \dots, \mathbf{v}_k]$ is an orthonormal basis to be determined. We can rewrite the above as

$$\mathbf{x}_i \approx \boldsymbol{\mu} + \mathbf{V} \boldsymbol{\alpha}_i,$$

where

$$\boldsymbol{\alpha}_i = \begin{bmatrix} \alpha_{i1} \\ \alpha_{i2} \\ \vdots \\ \alpha_{ik} \end{bmatrix},$$

with \mathbf{V} as a $n \times k$ matrix satisfying $\mathbf{V}^* \mathbf{V} = \mathbf{I}_k$. Now, we try to solve the optimization problem

$$\min_{\boldsymbol{\mu}, \mathbf{V}, \boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_n} I := \sum_{i=1}^n \|\mathbf{x}_i - (\boldsymbol{\mu} + \mathbf{V} \boldsymbol{\alpha}_i)\|_2^2.$$

Thus, we try to minimize the ℓ_2 -error across all vectors \mathbf{x}_i . (Unlike in the JL approach we do not strive for minimizing the error uniformly (within an ε -range) across all \mathbf{x}_i , but rather the *average error*.)

2.3 Solving the optimization problem

Fortunately we can separate this problem and first optimize over $\boldsymbol{\mu}$, then $\boldsymbol{\alpha}$, then over \mathbf{V} . (There are optimization problems which look similar but where you can't do this strategy of separation of variables.)

Let us first optimize with respect to $\boldsymbol{\mu}$. Without loss of generality, we can assume that $\sum_{i=1}^n \boldsymbol{\alpha}_i = 0$, because otherwise we could absorb the nonzero $\sum_i \boldsymbol{\alpha}_i$ into $\boldsymbol{\mu}$. Then:

$$\frac{\partial I}{\partial \boldsymbol{\mu}} = -2 \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{V} \boldsymbol{\alpha}_i).$$

Setting the right-hand side equal to zero, we get

$$\begin{aligned}\boldsymbol{\mu} &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \\ &= \boldsymbol{\mu}_n.\end{aligned}$$

Now let's optimize in α . We calculate:

$$\frac{\partial I}{\partial \alpha_i} = (\mathbf{x}_i - \boldsymbol{\mu} - \mathbf{V} \alpha_i)^* \mathbf{V}.$$

Setting the right-hand side equal to zero, we get

$$\alpha_i = \mathbf{V}^* (\mathbf{x}_i - \boldsymbol{\mu}).$$

Plugging in the expressions for μ and α_i into I , we get

$$I = \sum_{i=1}^n \|\mathbf{x}_i - \boldsymbol{\mu}_n - \mathbf{V} \mathbf{V}^* (\mathbf{x}_i - \boldsymbol{\mu}_n)\|_2^2$$

where $\mathbf{V} \mathbf{V}^*$ is an orthogonal projection matrix. Thus, letting $\mathbf{y}_i := \mathbf{x}_i - \boldsymbol{\mu}_n$,

$$I = \sum_{i=1}^n \|\mathbf{y}_i - \mathbf{V} \mathbf{V}^* \mathbf{y}_i\|_2^2.$$

Denote $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n]$. Then

$$\begin{aligned}\min_{\mathbf{V}: \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} \sum_{i=1}^n \|\mathbf{y}_i - \mathbf{V} \mathbf{V}^* \mathbf{y}_i\|_2^2 &= \min_{\mathbf{V}: \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} \text{trace} [(\mathbf{Y} - \mathbf{V} \mathbf{V}^* \mathbf{Y})^* (\mathbf{Y} - \mathbf{V} \mathbf{V}^* \mathbf{Y})] \\ &= \min_{\mathbf{V}: \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} \text{trace} [\mathbf{Y}^* (\mathbf{I} - \mathbf{V} \mathbf{V}^*) (\mathbf{I} - \mathbf{V} \mathbf{V}^*) \mathbf{Y}].\end{aligned}$$

Using properties of the trace (i.e., the *circular shift* property and linearity), and the fact that $(\mathbf{I} - \mathbf{V} \mathbf{V}^*)(\mathbf{I} - \mathbf{V} \mathbf{V}^*) = \mathbf{I} - \mathbf{V} \mathbf{V}^*$,

$$\begin{aligned}\min_{\mathbf{V}: \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} \sum_{i=1}^n \|\mathbf{y}_i - \mathbf{V} \mathbf{V}^* \mathbf{y}_i\|_2^2 &= \min_{\mathbf{V}: \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} \text{trace} [\mathbf{Y} \mathbf{Y}^* (\mathbf{I} - \mathbf{V} \mathbf{V}^*) (\mathbf{I} - \mathbf{V} \mathbf{V}^*)] \\ &= \min_{\mathbf{V}: \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} \text{trace} [\mathbf{Y} \mathbf{Y}^* (\mathbf{I} - \mathbf{V} \mathbf{V}^*)] \\ &= \min_{\mathbf{V}: \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} [\text{trace} (\mathbf{Y} \mathbf{Y}^*) - \text{trace} (\mathbf{Y} \mathbf{Y}^* \mathbf{V} \mathbf{V}^*)] \\ &= \min_{\mathbf{V}: \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} [\text{trace} (\mathbf{Y} \mathbf{Y}^*) - \text{trace} (\mathbf{Y} \mathbf{Y}^* \mathbf{V} \mathbf{V}^*)] \\ &= \min_{\mathbf{V}: \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} [\text{trace} (\mathbf{Y} \mathbf{Y}^*) - \text{trace} (\mathbf{V}^* \mathbf{Y} \mathbf{Y}^* \mathbf{V})]. \quad (1)\end{aligned}$$

But \mathbf{Y} does not depend on \mathbf{V} ! Hence, the minimum in (1) is independent of the expression $\text{trace}(\mathbf{Y}\mathbf{Y}^*)$ and thus coincides with the solution to

$$\max_{\mathbf{V} : \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} \frac{1}{n} \text{trace}(\mathbf{V}^* \mathbf{Y} \mathbf{Y}^* \mathbf{V}) = \max_{\mathbf{V} : \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} \text{trace}(\mathbf{V}^* \boldsymbol{\Sigma}_n \mathbf{V}).$$

Let $\boldsymbol{\Sigma}_n$ have eigenvalue decomposition

$$\boldsymbol{\Sigma}_n = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^*,$$

where $\lambda_i \geq 0$. (λ_i can't be negative because $\boldsymbol{\Sigma}_n$ is positive semidefinite.) The λ_i are the eigenvalues and the \mathbf{v}_i are the eigenvectors of $\boldsymbol{\Sigma}_n$. Since $\boldsymbol{\Sigma}_n$ is symmetric, the eigenvectors \mathbf{v}_i are orthogonal.

From linear algebra, we know that

$$\max_{\mathbf{V} : \mathbf{V}^* \mathbf{V} = \mathbf{I}_k} \text{trace}(\mathbf{V}^* \boldsymbol{\Sigma}_n \mathbf{V}) = \sum_{i=1}^k \lambda_i,$$

and *moreover*, the maximizer \mathbf{V} is the one given by $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_k]$. Hence, these particular \mathbf{v}_j give us the desired optimal orthonormal basis for our data \mathbf{x}_i .

2.4 Intuition for PCA

PCA takes the eigenvector decomposition of $\boldsymbol{\Sigma}_n$ and analyzes the projection of the centered data points (“centered” = subtract sample mean $\boldsymbol{\mu}_n$) on the k top eigenvectors of the sample covariance matrix $\boldsymbol{\Sigma}_n$ as *principal components*. (“ k top eigenvectors” = the eigenvectors associated with the largest k eigenvalues.)

2.5 Cost for PCA

The cost of this PCA procedure without using SVD is as follows. We need $\mathcal{O}(nd^2)$ operations to construct $\boldsymbol{\Sigma}_n$, and if you do the eigenvector decomposition in a traditional, naive sense, you need $\mathcal{O}(d^3)$ operations to find \mathbf{V} . However, the cost is a little bit cheaper via SVD, which we explain below.

Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ and

$$\mathbf{1}_n := \left\{ \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \right\} \text{ (n times).}$$

Then $\Sigma_n = \frac{1}{n} (\mathbf{X} - \mu_n \mathbf{1}_n^*) (\mathbf{X} - \mu_n \mathbf{1}_n^*)^*$. Now, the idea for saving computational cost is to just compute the SVD of $\mathbf{X} - \mu_n \mathbf{1}_n^*$.

The left singular vectors of $\mathbf{A} := \mathbf{X} - \mu_n \mathbf{1}_n^*$ are the same as the eigenvectors of $\mathbf{A}\mathbf{A}^* = \Sigma_n$, i.e., they are $\mathbf{v}_1, \dots, \mathbf{v}_n$. Therefore, the new cost via SVD is $\mathcal{O}(\min\{n^2d, nd^2\})$.

Now, from the *full* SVD we get *all* of $\mathbf{v}_1, \dots, \mathbf{v}_n$. But we really only want $\mathbf{v}_1, \dots, \mathbf{v}_k$, with $k < n$ (or even $k \ll n$). Computing only the top k singular vectors can be done in $\mathcal{O}(dnk)$ operations. In MATLAB, we can do this with the command `svds`¹. This is much faster than computing the full SVD. This computation of only the top k singular vectors is done via *Lanczos-type methods*².

We also note that *randomized SVD algorithms* can reduce this cost further to $\mathcal{O}(nd \log(k) + (n+d)k^2)$.

2.6 Another optimality property of the SVD

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m \geq n$ and let $\mathbf{A} = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{w}_i^*$. Denote $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{w}_i^*$ for $k < n$.

Given \mathbf{A} , then for any matrix \mathbf{B} of rank at most k , we have the following *best approximation* result:

$$\|\mathbf{A} - \mathbf{A}_k\|_{\text{op}} \leq \|\mathbf{A} - \mathbf{B}\|_{\text{op}},$$

and $\|\mathbf{A} - \mathbf{A}_k\|_{\text{op}} = \sigma_{k+1}$. In other words, \mathbf{A}_k is the *best rank- k approximation to \mathbf{A}* .

References

- [1] G.H. Golub and C.F. van Loan. *Matrix Computations*. Johns Hopkins, Baltimore, third edition, 1996.