

Theory related to PCA and SVD

Felix Held

Mathematical Sciences, Chalmers University of Technology and
University of Gothenburg

2019-04-08

Contents

1 Rayleigh quotient	1
1.1 The basic Rayleigh quotient	1
1.2 A more general Rayleigh quotient	2
2 Principal component analysis (PCA)	4
3 Singular value decomposition (SVD)	7
3.1 SVD and PCA	8
3.2 SVD and dimension reduction	9
3.3 SVD and orthogonal components	10
3.4 SVD and regression	11

1 Rayleigh quotient

1.1 The basic Rayleigh quotient

The Rayleigh quotient for a symmetric matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$ is a useful computational tool. It is defined for vectors $\mathbf{0} \neq \mathbf{x} \in \mathbb{R}^p$ as

$$J(\mathbf{x}) = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}. \quad (1.1)$$

Note that it is enough to normalize \mathbf{x} and calculate the Rayleigh quotient for the normalized vector since

$$J(\mathbf{x}) = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} = \frac{\|\mathbf{x}\|^2 \mathbf{x}^T \mathbf{A} \mathbf{x}}{\|\mathbf{x}\|^2 \mathbf{x}^T \mathbf{x}} = \frac{\frac{\mathbf{x}^T}{\|\mathbf{x}\|} \mathbf{A} \frac{\mathbf{x}}{\|\mathbf{x}\|}}{\frac{\mathbf{x}^T}{\|\mathbf{x}\|} \frac{\mathbf{x}}{\|\mathbf{x}\|}} = \frac{\mathbf{x}^T}{\|\mathbf{x}\|} \mathbf{A} \frac{\mathbf{x}}{\|\mathbf{x}\|}. \quad (1.2)$$

One says the Rayleigh quotient is scale-invariant. A common task is to find the $\hat{\mathbf{x}}$ that maximizes the Rayleigh quotient, i.e.

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}. \quad (1.3)$$

As noted above the Rayleigh quotient is scale-invariant. Therefore, for any solution $\hat{\mathbf{x}}$ of the maximization problem in Eq. (1.3) the vector $c \cdot \hat{\mathbf{x}}$ for arbitrary $c \in \mathbb{R}$ is a solution as well. To restrict the space of possible solutions we require that $\|\hat{\mathbf{x}}\| = 1$. This results in the optimization problem

$$\max_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} \quad \text{subject to} \quad \mathbf{x}^T \mathbf{x} = 1. \quad (1.4)$$

The Lagrangian (optimization course) of this problem is

$$L(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} - \lambda(\mathbf{x}^T \mathbf{x} - 1) \quad (1.5)$$

where λ is a Lagrange multiplier. We want to find a maximum of $L(\mathbf{x})$ and λ will be determined along the way. The gradient of $L(\mathbf{x})$ with respect to \mathbf{x} is

$$\frac{\partial L(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{A} \mathbf{x} - \lambda \mathbf{x} \quad (1.6)$$

and setting Eq. (1.6) equal to $\mathbf{0}$ leads to

$$\mathbf{A} \mathbf{x} = \lambda \mathbf{x}. \quad (1.7)$$

Since $\mathbf{x} \neq \mathbf{0}$ this is an eigenvalue equation and λ has to be one of p eigenvalues of \mathbf{A} . Using Eq. (1.7) in our original optimization problem Eq. (1.4) gives

$$\max_{\mathbf{x}, \|\mathbf{x}\|=1} \mathbf{x}^T \mathbf{A} \mathbf{x} = \max_{\substack{\mathbf{x}, \|\mathbf{x}\|=1 \\ \mathbf{A} \mathbf{x} = \lambda \mathbf{x}}} \lambda \mathbf{x}^T \mathbf{x} = \max_{\substack{\mathbf{x}, \|\mathbf{x}\|=1 \\ \mathbf{A} \mathbf{x} = \lambda \mathbf{x}}} \lambda \quad (1.8)$$

The optimization problem is therefore solved by an eigenvector \mathbf{x} of \mathbf{A} with $\|\mathbf{x}\| = 1$ such that the corresponding eigenvalue λ is maximal among all eigenvalues of \mathbf{A} . Note that there are always two solutions to this problem. For every \mathbf{x} with $\|\mathbf{x}\| = 1$ maximizing the Rayleigh quotient, the flipped vector $-\mathbf{x}$ is a solution with $\|\mathbf{x}\| = 1$ as well.

Note that since \mathbf{A} is real and symmetric, a theorem from linear algebra guarantees the existence of p real eigenvalues.

1.2 A more general Rayleigh quotient

A variant of the Rayleigh quotient assumes that there are two symmetric matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{p \times p}$ where \mathbf{B} is positive definite¹. The Rayleigh quotient is then defined for $\mathbf{0} \neq \mathbf{x} \in \mathbb{R}^p$ as

$$J(\mathbf{x}) = \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{B} \mathbf{x}}. \quad (1.9)$$

¹i.e. all eigenvalues are positive, $\mathbf{x}^T \mathbf{B} \mathbf{x} > 0$ for all $\mathbf{x} \neq \mathbf{0}$ and \mathbf{B} is invertible. Therefore, the numerator of the Rayleigh quotient is non-zero and positive as long as $\mathbf{x} \neq \mathbf{0}$.

As before $J(\mathbf{x})$ is invariant to scaling of \mathbf{x} and to make the optimization problem uniquely solvable (up to inverting the direction of the solution over the origin) we need to fixate a scaling. From linear algebra we know that there exists an orthogonal matrix² $\mathbf{U} \in \mathbb{R}^{p \times p}$ and a diagonal matrix³ $\mathbf{D} \in \mathbb{R}^{p \times p}$ such that $\mathbf{B} = \mathbf{U}\mathbf{D}\mathbf{U}^T$. Define $\mathbf{D}^{1/2} = \text{diag}(\sqrt{D_{ii}}, i = 1, \dots, p)$ and $\mathbf{B}^{1/2} = \mathbf{U}\mathbf{D}^{1/2}\mathbf{U}^T$, then $\mathbf{B}^{1/2}\mathbf{B}^{1/2} = \mathbf{B}$ ⁴.

In the previous section, we required to fix the scaling

$$\|\mathbf{x}\| = \mathbf{x}^T \mathbf{x} = 1. \quad (1.10)$$

This is very convenient for the applications below and it led for the numerator of the Rayleigh quotient to be one which made the proof above simpler. However, this is not equally convenient here since we have $\mathbf{x}^T \mathbf{B}\mathbf{x}$ in the numerator. We therefore require

$$\|\mathbf{B}^{1/2}\mathbf{x}\| = \mathbf{x}^T \mathbf{B}\mathbf{x} = 1, \quad (1.11)$$

which leads to the optimization problem

$$\max_{\mathbf{x}} \mathbf{x}^T \mathbf{A}\mathbf{x} \quad \text{subject to} \quad \mathbf{x}^T \mathbf{B}\mathbf{x} = 1. \quad (1.12)$$

Note that this is equivalent to solving

$$\max_{\mathbf{x}} \frac{\mathbf{x}^T \mathbf{A}\mathbf{x}}{\mathbf{x}^T \mathbf{B}\mathbf{x}} \quad \text{subject to} \quad \mathbf{x}^T \mathbf{x} = 1 \quad (1.13)$$

since for every solution $\hat{\mathbf{x}}$ to Eq. (1.13) the vector $\mathbf{z} = \hat{\mathbf{x}}/\|\mathbf{B}^{1/2}\hat{\mathbf{x}}\|$ also maximizes the Rayleigh quotient with $\mathbf{z}^T \mathbf{B}\mathbf{z} = 1$ and is therefore a solution to Eq. (1.12). On the other hand, every solution $\hat{\mathbf{x}}$ to Eq. (1.12) is a solution to Eq. (1.13) by setting $\mathbf{z} = \hat{\mathbf{x}}/\|\hat{\mathbf{x}}\|$.

Using a Lagrangian, calculating its gradient and setting it to zero (as above for the basic Rayleigh quotient) leads to

$$\mathbf{A}\mathbf{x} = \lambda \mathbf{B}\mathbf{x} \quad (1.14)$$

This is called a generalized eigenvalue problem for the symmetric matrices \mathbf{A} and \mathbf{B} . Using $\mathbf{B}^{1/2}$ as above we get

$$(\mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2})(\mathbf{B}^{1/2}\mathbf{x}) = \lambda (\mathbf{B}^{1/2}\mathbf{x}) \quad (1.15)$$

Defining $\mathbf{w} = \mathbf{B}^{1/2}\mathbf{x}$ results in

$$(\mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2})\mathbf{w} = \lambda \mathbf{w} \quad (1.16)$$

² $\mathbf{U}\mathbf{U}^T = \mathbf{U}^T\mathbf{U} = \mathbf{I}_n$

³The eigenvalues of \mathbf{U} are on the diagonal.

⁴Note that there was a similar construct in the lecture. A inverse covariance matrix Σ^{-1} was written as $\Sigma^{-1} = \mathbf{R}\mathbf{D}^{-1}\mathbf{R}^T$. With the definition $\Sigma^{-1/2} := \mathbf{D}^{-1/2}\mathbf{R}^T$ it holds that $(\Sigma^{-1/2})^T \Sigma^{-1/2} = \Sigma^{-1}$. Note that there is a transpose and above we are simply multiplying the “square root” matrices.

which is an eigenvalue problem for the symmetric matrix $\mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}$. This matrix is guaranteed to have p real-valued eigenvalues and corresponding eigenvectors \mathbf{w}_i for $i = 1, \dots, p$. Note that the rank of this matrix is determined by the rank of \mathbf{A} since $\mathbf{B}^{-1/2}$ is invertible. The original generalized eigenvectors \mathbf{x} can be recovered through

$$\mathbf{x} = \mathbf{B}^{-1/2}\mathbf{w}. \quad (1.17)$$

By using Eq. (1.14) in the original optimization problem in Eq. (1.12), show that this more general variant of the Rayleigh quotient is maximized for the vector \mathbf{x} such that $\mathbf{B}^{-1/2}\mathbf{x}$ is an eigenvector of the matrix $\mathbf{B}^{-1/2}\mathbf{A}\mathbf{B}^{-1/2}$ corresponding to its largest eigenvalue λ .

2 Principal component analysis (PCA)

When we have quantitative data, a natural choice of coordinate system is one where the axes point in the directions of largest variance and are orthogonal to each other. It is also natural to sort these in descending order, since most information will be gained by observing the most variable direction. Assume we have a data matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ with rows \mathbf{x}_i^T . To determine the first principal component we are looking for a direction \mathbf{r}_1 (a unit vector, i.e. $\|\mathbf{r}_1\| = 1$) in which the variance of \mathbf{X} is maximal. Define $s_i = \mathbf{r}_1^T \mathbf{x}_i$, which is the coefficient of \mathbf{x}_i projected onto \mathbf{r}_1 . The variance in the direction of \mathbf{r}_1 is

$$\sum_{i=1}^n (s_i - \bar{s})^2 = \sum_{i=1}^n (\mathbf{r}_1^T (\mathbf{x}_i - \bar{\mathbf{x}}))^2 \quad (2.1)$$

where $\bar{\mathbf{x}}$ is the mean over all observations. We want to find \mathbf{r}_1 such that the variance in Eq. (2.1) becomes maximal. Note that

$$\begin{aligned} \sum_{i=1}^n (\mathbf{r}_1^T (\mathbf{x}_i - \bar{\mathbf{x}}))^2 &= \sum_{i=1}^n \mathbf{r}_1^T (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{r}_1 \\ &= \mathbf{r}_1^T \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{r}_1 \\ &= (n-1) \mathbf{r}_1^T \hat{\Sigma} \mathbf{r}_1 \end{aligned} \quad (2.2)$$

where $\hat{\Sigma}$ is the empirical covariance matrix of the data. Since $\hat{\Sigma}$ is a symmetric matrix and it is required that $\|\mathbf{r}_1\| = 1$, maximizing Eq. (2.2) is equivalent to solving the basic Rayleigh quotient maximisation problem in Section 1. It therefore follows that \mathbf{r}_1 is an eigenvector of $\hat{\Sigma}$ corresponding to its largest eigenvalue λ_1 . Since we required \mathbf{r}_1 to be of length one, this problem is solved uniquely up to sign (i.e. $-\mathbf{r}_1$ is also a solution). Note especially that the variance of the s_i , that we tried to maximize in the original problem (Eq. (2.1)) is equal to λ_1 .

Assume we have found the first $m-1 < p$ principal components $\mathbf{r}_1, \dots, \mathbf{r}_{m-1}$ corresponding to the eigenvalues $\lambda_1 \geq \dots \geq \lambda_{m-1}$ of $\hat{\Sigma}$. From linear algebra we know that

a square matrix \mathbf{P} is an orthogonal projection matrix if

$$\mathbf{P}^2 = \mathbf{P} = \mathbf{P}^T. \quad (2.3)$$

Projecting a vector onto $\mathbf{r}_1, \dots, \mathbf{r}_{m-1}$ is accomplished by the orthogonal projection matrix

$$\mathbf{P} = \sum_{i=1}^{m-1} \mathbf{r}_i \mathbf{r}_i^T \quad (2.4)$$

Another property of orthogonal projection matrices is that $\mathbf{I} - \mathbf{P}$ is an orthogonal projection matrix onto the orthogonal complement of the subspace that \mathbf{P} was projecting on, i.e. $\mathbf{I} - \mathbf{P}_1$ is a projection matrix onto the space of vectors which are orthogonal to $\mathbf{r}_1, \dots, \mathbf{r}_{m-1}$.

Project the data into this space, after having found the first $m - 1$ principal components, i.e. define

$$\mathbf{X}_{m-1} = \mathbf{X} \left(\mathbf{I}_p - \sum_{i=1}^{m-1} \mathbf{r}_i \mathbf{r}_i^T \right) \quad (2.5)$$

and

$$\hat{\Sigma}_{m-1} = \frac{\mathbf{X}_{m-1}^T \mathbf{X}_{m-1}}{n-1} = \left(\mathbf{I}_p - \sum_{i=1}^{m-1} \mathbf{r}_i \mathbf{r}_i^T \right)^T \hat{\Sigma} \left(\mathbf{I}_p - \sum_{i=1}^{m-1} \mathbf{r}_i \mathbf{r}_i^T \right). \quad (2.6)$$

The new data matrix is constant (no variance) along the directions of $\mathbf{r}_1, \dots, \mathbf{r}_{m-1}$. Finding the most variable direction now means to solve

$$\max_{\mathbf{r}} \mathbf{r}^T \hat{\Sigma}_{m-1} \mathbf{r} \quad \text{subject to} \quad \mathbf{r}^T \mathbf{r} = 1. \quad (2.7)$$

This is solved by an eigenvector \mathbf{r}_m of $\hat{\Sigma}_{m-1}$ corresponding to its largest eigenvalue λ_m , i.e. $\hat{\Sigma}_{m-1} \mathbf{r}_m = \lambda_m \mathbf{r}_m$. It turns out that for $i = 1, \dots, m-1$

$$\begin{aligned} \mathbf{r}_m^T \mathbf{r}_i &= \frac{1}{\lambda_m} \mathbf{r}_m^T \hat{\Sigma}_{m-1} \mathbf{r}_i \\ &= \frac{1}{\lambda_m} \mathbf{r}_m^T \left(\mathbf{I}_p - \sum_{j=1}^{m-1} \mathbf{r}_j \mathbf{r}_j^T \right)^T \underbrace{\hat{\Sigma} \left(\mathbf{I}_p - \sum_{j=1}^{m-1} \mathbf{r}_j \mathbf{r}_j^T \right) \mathbf{r}_i}_{=0} = 0 \end{aligned} \quad (2.8)$$

as well as

$$\begin{aligned}
\lambda_m \mathbf{r}_m &= \hat{\Sigma}_{m-1} \mathbf{r}_m \\
&= \left(\mathbf{I}_p - \sum_{j=1}^{m-1} \mathbf{r}_j \mathbf{r}_j^T \right)^T \hat{\Sigma} \left(\mathbf{I}_p - \sum_{j=1}^{m-1} \mathbf{r}_j \mathbf{r}_j^T \right) \mathbf{r}_m \\
&= \left(\mathbf{I}_p - \sum_{j=1}^{m-1} \mathbf{r}_j \mathbf{r}_j^T \right) \hat{\Sigma} \mathbf{r}_m \\
&= \hat{\Sigma} \mathbf{r}_m - \sum_{j=1}^{m-1} \mathbf{r}_j \underbrace{\mathbf{r}_j^T \hat{\Sigma} \mathbf{r}_m}_{=\lambda_j \mathbf{r}_j^T \mathbf{r}_m = 0} \\
&= \hat{\Sigma} \mathbf{r}_m
\end{aligned} \tag{2.9}$$

which shows that \mathbf{r}_m is orthogonal to $\mathbf{r}_1, \dots, \mathbf{r}_{m-1}$ and an eigenvector of $\hat{\Sigma}$. In addition, since $\lambda_1 \geq \dots \geq \lambda_{m-1}$ are the $m-1$ largest eigenvalues of $\hat{\Sigma}$, it must hold that $\lambda_{m-1} \geq \lambda_m$.

As a procedural way of calculating the PCA of a data matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$, one has to follow the following steps

1. Centre and standardize the columns of the data matrix \mathbf{X} (the variables)
2. Calculate the empirical covariance matrix

$$\hat{\Sigma} = \frac{1}{n-1} \mathbf{X}^T \mathbf{X} \tag{2.10}$$

3. Determine the eigenvalues λ_i for $i = 1, \dots, p$ of $\hat{\Sigma}$ and a set of p corresponding orthonormal eigenvectors \mathbf{r}_i such that

$$\hat{\Sigma} \mathbf{r}_i = \lambda_i \mathbf{r}_i, \|\mathbf{r}_i\| = 1, i = 1, \dots, p \quad \text{and} \quad \mathbf{r}_i^T \mathbf{r}_j = 0, i \neq j \tag{2.11}$$

as well as

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \tag{2.12}$$

4. Set

$$\mathbf{R} = (\mathbf{r}_1, \dots, \mathbf{r}_p) \in \mathbb{R}^{p \times p} \quad \text{and} \quad \mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_p) \tag{2.13}$$

so that

$$\hat{\Sigma} = \mathbf{R} \mathbf{D} \mathbf{R}^T \tag{2.14}$$

5. The vectors \mathbf{r}_i are the principal component directions, the projections $\mathbf{r}_i^T \mathbf{x}$ are called principal components and the corresponding eigenvalues λ_i are the variance of the data in the direction of the principal component.

PCA can be used to reduce the dimension of the data. Since the principal components account for less variance in every step, it is possible that there is little information in the last principal components. An important result from linear algebra is that for a matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$ with eigenvalues μ_1, \dots, μ_p it holds that

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^p \mu_i. \quad (2.15)$$

The empirical covariance matrix has the variance of each variable on its diagonal and therefore

$$\text{tr}(\hat{\Sigma}) = \sum_{i=1}^p s^2(\mathbf{X}_{\cdot i}) = \sum_{i=1}^p \lambda_i, \quad (2.16)$$

where $\mathbf{X}_{\cdot i}$ is the i -th column of the data matrix and s^2 is the empirical variance. If the variables are standardized then $s^2(\mathbf{X}_{\cdot i}) = 1$ for all i and therefore $\text{tr}(\hat{\Sigma}) = p$. This means in particular that the mean of the eigenvalues will be 1⁵. A typical criterion for considering a principal component as important is that the corresponding eigenvalue is larger than the mean of all eigenvalues (in case of standardised data: larger than one). A tool for the analysis of the information behind the principal components is a *scree plot*. In a scree plot, the eigenvalues (variances) are plotted as a function of the index of the principal components. It is a way to quickly see how many principal components are of interest.

3 Singular value decomposition (SVD)

The singular value decomposition (SVD) of a matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$, $n \geq p$, splits the data matrix into a product of three matrices

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T \quad (3.1)$$

where $\mathbf{U} \in \mathbb{R}^{n \times p}$ has orthonormal columns, $\mathbf{D} \in \mathbb{R}^{p \times p}$ is a diagonal matrix, and $\mathbf{V} \in \mathbb{R}^{p \times p}$ is an orthogonal matrix. Note that

$$\mathbf{U}^T \mathbf{U} = \mathbf{I}_p \quad \text{and} \quad \mathbf{V}\mathbf{V}^T = \mathbf{V}^T \mathbf{V} = \mathbf{I}_p \quad (3.2)$$

Note that if $n > p$ it cannot hold that $\mathbf{U}\mathbf{U}^T = \mathbf{I}_n$ ⁶. The matrix \mathbf{D} is diagonal and contains the *singular values* d_i . These are typically sorted such that $d_{i+1} \leq d_i$.

The orthogonality properties in Eq. (3.2) can now be used to derive the following equations

$$\begin{aligned} \mathbf{X}\mathbf{X}^T \mathbf{U} &= \mathbf{U}\mathbf{D}\mathbf{V}^T \mathbf{V}\mathbf{D}\mathbf{U}^T \mathbf{U} = \mathbf{U}\mathbf{D}^2 \\ \mathbf{X}^T \mathbf{X}\mathbf{V} &= \mathbf{V}\mathbf{D}\mathbf{U}^T \mathbf{U}\mathbf{D}\mathbf{V}^T \mathbf{V} = \mathbf{V}\mathbf{D}^2 \end{aligned} \quad (3.3)$$

⁵ $\sum_{i=1}^p \lambda_i / p = 1$

⁶A set of maximally p vectors can be linearly independent in \mathbb{R}^p and the equation above would imply that $n > p$ vectors in \mathbb{R}^p are orthogonal and thus linearly independent.

Since \mathbf{D} is a diagonal matrix, this reduces the problem of determining \mathbf{U} and \mathbf{V} to solving a series of eigenvalue problems

$$\mathbf{X}\mathbf{X}^T\mathbf{u}_i = d_i^2\mathbf{u}_i \quad \text{and} \quad \mathbf{X}^T\mathbf{X}\mathbf{v}_i = d_i^2\mathbf{v}_i, \quad i = 1, \dots, p \quad (3.4)$$

Since $n \geq p$, the matrix $\mathbf{X}^T\mathbf{X} \in \mathbb{R}^{p \times p}$ is as large as or smaller than $\mathbf{X}\mathbf{X}^T \in \mathbb{R}^{n \times n}$. It is therefore more computationally effective to calculate \mathbf{V} by solving the p eigenvalue problems in Eq. (3.4) and then arrive at \mathbf{U} by projecting the observations in \mathbf{X} on the space spanned by the columns of \mathbf{V} and scaling them by the inverse of the singular values, i.e.

$$\mathbf{U} = \mathbf{X}\mathbf{V}\mathbf{D}^{-1}. \quad (3.5)$$

Note that this approach requires that there are no singular values equal to zero (which is allowed and possible in general).

For $n < p$ the SVD can still be calculated. Note that the SVD can be calculated as above for $\mathbf{X}^T \in \mathbb{R}^{p \times n}$. We get matrices $\mathbf{V} \in \mathbb{R}^{p \times n}$ (orthonormal columns), a diagonal matrix $\mathbf{D} \in \mathbb{R}^{n \times n}$, and $\mathbf{U} \in \mathbb{R}^{n \times n}$ (orthogonal matrix) such that

$$\mathbf{X}^T = \mathbf{V}\mathbf{D}\mathbf{U}^T. \quad (3.6)$$

By transposing again we get

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T. \quad (3.7)$$

Note that this time \mathbf{U} is square and \mathbf{V} is rectangular.

3.1 SVD and PCA

The SVD of a data matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ is usually used to calculate the principal components in the data. Assume that $n \geq p$ for now. Assume that the variables in \mathbf{X} have been centred and scaled. Note that the empirical covariance matrix of \mathbf{X} is

$$\hat{\Sigma} = \frac{\mathbf{X}^T\mathbf{X}}{n-1}. \quad (3.8)$$

Using the SVD of $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ this leads to

$$\hat{\Sigma} = \frac{1}{n-1}\mathbf{V}\mathbf{D}\mathbf{U}^T\mathbf{U}\mathbf{D}\mathbf{V}^T = \mathbf{V}\left(\frac{1}{n-1}\mathbf{D}^2\right)\mathbf{V}^T. \quad (3.9)$$

Comparing this with an orthogonal decomposition of a symmetric real matrix, we get that the eigenvalues of $\hat{\Sigma}$ are on the diagonal of the matrix $\mathbf{D}^2/(n-1)$, and the columns of \mathbf{V} are the corresponding eigenvectors. This determines the principal components and the corresponding variance explained by each component.

For $n < p$ the calculations above work out the same, however, \mathbf{V} is not a square matrix any longer. The interpretation of eigenvalues and principal components still holds up since

$$\hat{\Sigma}\mathbf{V} = \mathbf{V}\left(\frac{1}{n-1}\mathbf{D}^2\right) \quad (3.10)$$

but note that the eigenvalues $\lambda_{n+1}, \dots, \lambda_p$ are zero.

3.2 SVD and dimension reduction

When using SVD for dimension reduction, we want to reduce the number of variables. Therefore we want to find a subspace of \mathbb{R}^p of dimension $q \leq \min(n, p)$ such that the projections of the observations into this subspace are as similar to the original observations in the Euclidean norm as possible, i.e. the objective is to find a subspace $\hat{\mathcal{S}} \subset \mathbb{R}^p$, such that

$$\hat{\mathcal{S}} = \arg \min_{\mathcal{S}} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{P}_{\mathcal{S}} \mathbf{x}_i\|_2^2 \quad (3.11)$$

where $\mathbf{P}_{\mathcal{S}}$ is the orthogonal projection of \mathbb{R}^p onto the subspace \mathcal{S} . Note that

$$\begin{aligned} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{P}_{\mathcal{S}} \mathbf{x}_i\|_2^2 &= \sum_{i=1}^n (\mathbf{x}_i - \mathbf{P}_{\mathcal{S}} \mathbf{x}_i)^T (\mathbf{x}_i - \mathbf{P}_{\mathcal{S}} \mathbf{x}_i) \\ &= \sum_{i=1}^n [\mathbf{x}_i^T \mathbf{x}_i - 2\mathbf{x}_i^T \mathbf{P}_{\mathcal{S}} \mathbf{x}_i + \mathbf{x}_i^T \mathbf{P}_{\mathcal{S}}^T \mathbf{P}_{\mathcal{S}} \mathbf{x}_i] \\ &= \sum_{i=1}^n [\mathbf{x}_i^T \mathbf{x}_i - \mathbf{x}_i^T \mathbf{P}_{\mathcal{S}}^T \mathbf{P}_{\mathcal{S}} \mathbf{x}_i] \\ &= \sum_{i=1}^n [\|\mathbf{x}_i\|_2^2 - \|\mathbf{P}_{\mathcal{S}} \mathbf{x}_i\|_2^2] \end{aligned} \quad (3.12)$$

This shows that minimizing the squared distance between the observations and their projections on the subspace is equivalent to maximizing the squared length of the projected vectors⁷. Let's start by looking for the best one-dimensional subspace \mathcal{S}_1 . Every one-dimensional vector space is spanned by a vector \mathbf{r}_1 and we can assume that $\|\mathbf{r}_1\|_2 = 1$. The projection matrix onto a one-dimensional vector space is $\mathbf{P} = \mathbf{r}_1 \mathbf{r}_1^T$ and thus we are trying to maximize

$$\begin{aligned} \sum_{i=1}^n \|\mathbf{r}_1 \mathbf{r}_1^T \mathbf{x}_i\|_2^2 &= \sum_{i=1}^n (\mathbf{r}_1^T \mathbf{x}_i)^2 \\ &= \mathbf{r}_1^T \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{r}_1 \\ &= (n-1) \mathbf{r}_1^T \hat{\Sigma} \mathbf{r}_1 \end{aligned} \quad (3.13)$$

The term we are maximizing is the Rayleigh quotient and therefore the optimal subspace is spanned by \mathbf{r}_1 , the eigenvector of $\hat{\Sigma}$ with the largest corresponding eigenvalue. So the optimal one-dimensional sub-space that is closest to the data is spanned by the first principal component, which is also the first column of the matrix \mathbf{V} in the SVD of \mathbf{X} .

⁷Since $\sum_{i=1}^n \|\mathbf{x}_i\|_2^2$ is a constant given a dataset. The only object we can control is the subspace we project into and therefore we can only change $\sum_{i=1}^n \|\mathbf{P}_{\mathcal{S}} \mathbf{x}_i\|_2^2$.

Similar arguments as for PCA lead to the conclusion that the best q -dimensional subspace to approximate the data in is the space spanned by the first q principal component directions, which are also the first q columns in the matrix \mathbf{V} in the SVD of \mathbf{X} . Note that for $n < p$ the data is maximally n dimensional and any approximating subspace must therefore have dimension $< n$.

3.3 SVD and orthogonal components

Another interpretation of SVD is that it describes a method to describe the data as a structure of orthogonal components. Let $\mathbf{u}_i \in \mathbb{R}^n$ be the columns of \mathbf{U} and $\mathbf{v}_i \in \mathbb{R}^{\min(n,p)}$ the columns of \mathbf{V} . Then

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T = \sum_{i=1}^{\min(n,p)} d_i \mathbf{u}_i \mathbf{v}_i^T. \quad (3.14)$$

Each matrix $\mathbf{u}_i \mathbf{v}_i^T$ is of rank 1 and these matrices are scaled by the singular values d_i . As we have seen above in Section 3.2, the optimal $q \leq \min(n, p)$ subspace to approximate \mathbf{X} is spanned by the first q columns of \mathbf{V} and the projection on that subspace is

$$\mathbf{P}_q = \sum_{i=1}^q \mathbf{v}_i \mathbf{v}_i^T. \quad (3.15)$$

Projecting \mathbf{X} on this optimal q -dimensional subspace using the projection in Eq. (3.15) and the representation in Eq. (3.14) leads to

$$\mathbf{X}_q = \mathbf{X}\mathbf{P}_q = \left(\sum_{i=1}^{\min(n,p)} d_i \mathbf{u}_i \mathbf{v}_i^T \right) \left(\sum_{j=1}^q \mathbf{v}_j \mathbf{v}_j^T \right) = \sum_{i=1}^q d_i \mathbf{u}_i \mathbf{v}_i^T. \quad (3.16)$$

So projecting the data into this optimal q -dimensional subspace simply means to only keep components 1 to q . The approximation error in Frobenius norm⁸ is

$$\begin{aligned}
\|\mathbf{X} - \mathbf{X}_q\|_F^2 &= \left\| \sum_{i=q+1}^{\min(n,p)} d_i \mathbf{u}_i \mathbf{v}_i^T \right\|_F^2 \\
&= \text{tr} \left[\left(\sum_{i=q+1}^{\min(n,p)} d_i \mathbf{u}_i \mathbf{v}_i^T \right)^T \left(\sum_{j=q+1}^{\min(n,p)} d_j \mathbf{u}_j \mathbf{v}_j^T \right) \right] \\
&= \text{tr} \left[\sum_{i,j=q+1}^{\min(n,p)} d_i d_j \mathbf{v}_i \mathbf{u}_i^T \mathbf{u}_j \mathbf{v}_j^T \right] \\
&= \text{tr} \left[\sum_{i,j=q+1}^{\min(n,p)} d_i d_j \mathbf{v}_i \mathbf{v}_j^T \right] \tag{3.17} \\
&= \sum_{k=1}^{\min(n,p)} \sum_{i,j=q+1}^{\min(n,p)} d_i d_j v_{ik} v_{jk} \\
&= \sum_{i,j=q+1}^{\min(n,p)} d_i d_j \sum_{k=1}^{\min(n,p)} v_{ik} v_{jk} \\
&= \sum_{i,j=q+1}^{\min(n,p)} d_i d_j \mathbb{1}(i=j) = \sum_{i=q+1}^{\min(n,p)} d_i^2
\end{aligned}$$

3.4 SVD and regression

In addition to dimension reduction and easy determination of the optimal q -dimensional approximation to the data, SVD can also be a useful tool in regression.

Recall the linear regression problem with response vector $\mathbf{y} \in \mathbb{R}^n$ and design matrix $\mathbf{X} \in \mathbb{R}^{n \times (p+1)}$, where we adopt the convention that the first column of \mathbf{X} is a vector of 1's to encode the intercept. The variable y is modelled as

$$\mathbf{y} = \mathbf{x}^T \boldsymbol{\beta} + \varepsilon \tag{3.18}$$

where $\mathbf{x} \in \mathbb{R}^{p+1}$ is a vector of predictors with $x_1 = 1$ and $\boldsymbol{\beta} \in \mathbb{R}^{p+1}$ are the regression coefficients. The variable ε is the error and typically one assumes that

$$\varepsilon \sim \text{Normal}(0, \sigma^2) \tag{3.19}$$

for some (possibly unknown) variance σ^2 . A solution to the regression model can be found with least squares, i.e. solving

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^{p+1}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \tag{3.20}$$

⁸This is a matrix norm defined by $\|\mathbf{X}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^p x_{ij}^2$

The solution to the least squares problem is given by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad (3.21)$$

which requires $n \geq p$ for the inversion to be possible. The SVD of $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$ can be plugged into this equation to arrive at

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{U} \mathbf{D} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{D}^2 \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} \mathbf{D}^{-2} \mathbf{V}^T \mathbf{V} \mathbf{D} \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} \mathbf{D}^{-1} \mathbf{U}^T \mathbf{y} \end{aligned} \quad (3.22)$$

The expression

$$\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{V} \mathbf{D}^{-1} \mathbf{U}^T \quad (3.23)$$

is called the Moore-Penrose pseudo-inverse of \mathbf{X} . While it is therefore possible to obtain a least squares solution through SVD, there are other simpler algorithms (e.g. QR decomposition) which are preferable.

The fitted values for \mathbf{y} are then

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{U} \mathbf{D} \mathbf{V}^T \mathbf{V} \mathbf{D}^{-1} \mathbf{U}^T \mathbf{y} = \mathbf{U} \mathbf{U}^T \mathbf{y}. \quad (3.24)$$

Note that $\mathbf{U} \mathbf{U}^T$ is an orthogonal projection matrix, projecting \mathbf{y} onto the column space of \mathbf{U} .

For ridge regression the problem to be solved is

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X} \boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2 \quad (3.25)$$

for some $\lambda \geq 0$. Assume that the response \mathbf{y} and the columns of \mathbf{X} are centred. The solution is given by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}. \quad (3.26)$$

Using the SVD of $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$ leads to

$$\hat{\boldsymbol{\beta}} = \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} \mathbf{D} \mathbf{U}^T \mathbf{y} = \sum_{i=1}^p \frac{d_i}{d_i^2 + \lambda} \mathbf{v}_i \mathbf{u}_i^T \mathbf{y}. \quad (3.27)$$

It can be seen that λ can lead to stability in the calculation of the fractions $d_i/(d_i^2 + \lambda)$. If d_i is small, the fraction would become big if $\lambda = 0$. Increasing lambda decreases the magnitude of the fractions and increases numerical stability. Also, an increase in λ decreases the influence of each term in the same and therefore shrinks coefficients towards 0.