

Some classical matrix factorisation techniques

Felix Held¹

¹Mathematical Sciences, University of Gothenburg and Chalmers University of Technology,
Gothenburg, Sweden, felix.held@chalmers.se

Started: 2020-01-10, Last update: 2020-04-02

Contents

1	Introduction	1
2	Rayleigh quotient	2
2.1	The basic Rayleigh quotient	2
2.2	The generalised Rayleigh quotient	4
3	Principal component analysis (PCA)	5
4	Singular value decomposition (SVD)	8
4.1	SVD and PCA	9
4.2	SVD and dimension reduction	9
4.3	SVD and orthogonal components	10
4.4	SVD and regression	11
5	Factor Analysis (FA)	13
6	Probabilistic Principal Component Analysis (PPCA)	14
7	Canonical Correlation Analysis (CCA)	15
8	Independent Component Analysis (ICA)	17

1 Introduction

The purpose of this document is to give an overview over classical/traditional matrix factorisation techniques and how they are interpreted. For a start I want to focus on the three classical techniques Principal Component Analysis (PCA), Independent Component Analysis (ICA), and Canonical Correlation Analysis (CCA). As a working tool,

the Rayleigh quotient and the Singular Value Decomposition (SVD) of a matrix will be introduced as well.

The exposition assumes that the reader has a background in multivariate calculus and basic multivariate statistics. The material presented here is somewhat standard and presented in many books and other material across the internet. Some of my sources were [Falk et al. \(2002\)](#), [Härdle and Simar \(2007\)](#) and [Murphy \(2012\)](#).

2 Rayleigh quotient

2.1 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

$$J(\mathbf{x}) = \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}. \quad (1)$$

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

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

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

$$\hat{\mathbf{x}} \in \arg \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}. \quad (3)$$

As noted above the Rayleigh quotient is scale-invariant. Therefore, for any solution $\hat{\mathbf{x}}$ of the maximization problem in [Eq. \(3\)](#) the vector $c \cdot \hat{\mathbf{x}}$ for arbitrary $c \in \mathbb{R} \setminus \{0\}$ is a solution as well. To restrict the space of possible solutions and to make the optimization problem identifiable we require that $\|\hat{\mathbf{x}}\| = 1$. This results in the equality-constrained minimization problem

$$\min_{\mathbf{x} \neq \mathbf{0}} -\mathbf{x}^\top \mathbf{A} \mathbf{x} \quad \text{subject to} \quad \mathbf{x}^\top \mathbf{x} = 1. \quad (4)$$

The Lagrangian ([Boyd and Vandenberghe, 2004](#), Section 5.1) of this constrained maximization problem is

$$L(\mathbf{x}, \lambda) = -\mathbf{x}^\top \mathbf{A} \mathbf{x} + \lambda(\mathbf{x}^\top \mathbf{x} - 1) \quad (5)$$

where λ is a Lagrange multiplier. The Lagrange dual function is therefore

$$g(\lambda) = \inf_{\mathbf{x} \neq \mathbf{0}} \mathbf{x}^\top (-\mathbf{A} + \lambda \mathbf{I}_p) \mathbf{x} - \lambda \quad (6)$$

For the infimum to be finite $-\mathbf{A} + \lambda \mathbf{I}_p$ has to be positive semi-definite (psd). If $-\mathbf{A} + \lambda \mathbf{I}_p$ has at least one negative eigenvalue μ with corresponding eigenvector \mathbf{u} , then

$c^2 \mathbf{u}^\top (-\mathbf{A} + \lambda \mathbf{I}_p) \mathbf{u} = c^2 \mu \|\mathbf{u}\|_2^2 \rightarrow -\infty$ for $c \rightarrow \infty$. This implies that the infimum is $-\infty$ in the non-psd case.

Since \mathbf{A} is symmetric, let $\mathbf{U} \in \mathbb{R}^{p \times p}$ contain an orthonormal basis of eigenvectors of \mathbf{A} , with the i -th basis vector \mathbf{u}_i in the i -th row of \mathbf{U} , and \mathbf{D} the corresponding matrix of eigenvalues of \mathbf{A} . Then $\mathbf{A} = \mathbf{U}^\top \mathbf{D} \mathbf{U}$ and

$$\mathbf{x}^\top (-\mathbf{A} + \lambda \mathbf{I}_p) \mathbf{x} = \mathbf{x}^\top \mathbf{U}^\top (-\mathbf{D} + \lambda \mathbf{I}_p) \mathbf{U} \mathbf{x} = \sum_{i=1}^p (\lambda - \mu_i) (\mathbf{u}_i^\top \mathbf{x})^2. \quad (7)$$

This implies especially that the matrix $-\mathbf{A} + \lambda \mathbf{I}_p$ is psd if and only if $\lambda \geq \mu_{\max}(\mathbf{A})$ where $\mu_{\max}(\mathbf{A})$ is the largest eigenvalue of \mathbf{A} . Since every $\mathbf{x} = \sum_{i=1}^p a_i \mathbf{u}_i$ for some a_i it follows that $\sum_{i=1}^p (\lambda - \mu_i) (\mathbf{u}_i^\top c \mathbf{x})^2 = c^2 \sum_{i=1}^p (\lambda - \mu_i) a_i^2 \rightarrow 0$ for $c \rightarrow 0$. This implies that the infimum is $-\lambda$ for $\lambda \geq \mu_{\max}(\mathbf{A})$ in the psd case.

The Lagrange dual function for the minimization problem in Eq. (4) is therefore given by

$$g(\lambda) = \begin{cases} -\lambda & \text{if } \lambda \geq \mu_{\max}(\mathbf{A}) \\ -\infty & \text{otherwise} \end{cases} \quad (8)$$

which implies the dual problem

$$\sup_{\lambda} -\lambda \quad \text{subject to} \quad \lambda \geq \mu_{\max}(\mathbf{A}). \quad (9)$$

Clearly, this optimization problem is solved for $\hat{\lambda} = \mu_{\max}(\mathbf{A})$.

It now follows that, since the optimal function value $d^* = -\mu_{\max}(\mathbf{A})$ of the dual problem is a lower bound on the optimal function value p^* of the primal problem (Eq. (4)), we get by choosing \mathbf{u} as an unit eigenvector corresponding to $\mu_{\max}(\mathbf{A})$

$$-\mu_{\max}(\mathbf{A}) \leq p^* \leq -\mathbf{u}^\top \mathbf{A} \mathbf{u} = -\mu_{\max}(\mathbf{A}). \quad (10)$$

This shows that $\mathbf{u} \in \arg \max_{\mathbf{x} \neq \mathbf{0}} J(\mathbf{x})$. Note that $-\mathbf{u}$ is also a solution. In case the eigenvalue $\mu_{\max}(\mathbf{A})$ appears repeatedly, any \mathbf{u} in its eigenspace with length $\|\mathbf{u}\| = 1$ will be a solution of the primal problem in Eq. (4).

A more hand-waving derivation is as follows: We want to find a maximum of $L(\mathbf{x}, \lambda)$ with respect to \mathbf{x} and λ will be determined along the way. The gradient of $L(\mathbf{x}, \lambda)$ with respect to \mathbf{x} is

$$\frac{\partial L(\mathbf{x}, \lambda)}{\partial \mathbf{x}} = -\mathbf{A} \mathbf{x} + \lambda \mathbf{x} \quad (11)$$

and setting Eq. (11) equal to $\mathbf{0}$ leads to

$$\mathbf{A} \mathbf{x} = \lambda \mathbf{x}. \quad (12)$$

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. (12) in our original optimization problem Eq. (4) gives

$$\max_{\mathbf{x}, \|\mathbf{x}\|=1} \mathbf{x}^\top \mathbf{A} \mathbf{x} = \max_{\substack{\mathbf{x}, \|\mathbf{x}\|=1 \\ \mathbf{A} \mathbf{x} = \lambda \mathbf{x}}} \lambda \mathbf{x}^\top \mathbf{x} = \max_{\substack{\mathbf{x}, \|\mathbf{x}\|=1 \\ \mathbf{A} \mathbf{x} = \lambda \mathbf{x}}} \lambda \quad (13)$$

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 sign-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.

2.2 The generalised 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

$$J(\mathbf{x}) = \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{B} \mathbf{x}}. \quad (14)$$

As above, $J(\mathbf{x})$ 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 matrix $\mathbf{D} \in \mathbb{R}^{p \times p}$ such that $\mathbf{B} = \mathbf{U} \mathbf{D} \mathbf{U}^\top$. Define $\mathbf{D}^{1/2} := \text{diag}(\sqrt{D_{ii}}, i = 1, \dots, p)$ and $\mathbf{B}^{1/2} := \mathbf{U} \mathbf{D}^{1/2} \mathbf{U}^\top$, then $\mathbf{B}^{1/2} \mathbf{B}^{1/2} = \mathbf{B}$ ³. By convention, the inverse of $\mathbf{B}^{1/2}$ is denoted as $\mathbf{B}^{-1/2}$.

To fixate a scaling for the input to the generalised Rayleigh quotient we require

$$\|\mathbf{B}^{1/2} \mathbf{x}\| = \mathbf{x}^\top \mathbf{B} \mathbf{x} = 1, \quad (15)$$

which leads to the optimization problem⁴

$$\max_{\mathbf{x}} \mathbf{x}^\top \mathbf{A} \mathbf{x} \quad \text{subject to} \quad \mathbf{x}^\top \mathbf{B} \mathbf{x} = 1. \quad (17)$$

Since $\mathbf{B}^{1/2}$ is of full-rank⁵, we can make the change of variable $\mathbf{w} \leftarrow \mathbf{B}^{1/2} \mathbf{x}$ in Eq. (17), resulting in

$$\max_{\mathbf{w}} \mathbf{w}^\top \mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2} \mathbf{w} \quad \text{subject to} \quad \mathbf{w}^\top \mathbf{w} = 1. \quad (18)$$

This is the optimization problem for the basic Rayleigh quotient with symmetric matrix $\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2}$. Any $\hat{\mathbf{w}}$ which is a unit eigenvector corresponding to $\mu_{\max}(\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2})$

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

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

³ $\mathbf{B}^{1/2}$ is sometimes called the matrix square root of \mathbf{B} .

⁴Note that this is equivalent to solving

$$\max_{\mathbf{x}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{B} \mathbf{x}} \quad \text{subject to} \quad \mathbf{x}^\top \mathbf{x} = 1 \quad (16)$$

since for every solution $\hat{\mathbf{x}}$ to Eq. (16) the vector $\mathbf{z} = \hat{\mathbf{x}} / \|\mathbf{B}^{1/2} \hat{\mathbf{x}}\|$ also maximizes the Rayleigh quotient with $\mathbf{z}^\top \mathbf{B} \mathbf{z} = 1$ and is therefore a solution to Eq. (17). On the other hand, every solution $\hat{\mathbf{x}}$ to Eq. (17) induces a solution to Eq. (16) by setting $\hat{\mathbf{z}} = \hat{\mathbf{x}} / \|\hat{\mathbf{x}}\|$.

⁵Follows from \mathbf{B} being psd.

is therefore a solution to [Eq. \(18\)](#). Corresponding solutions for [Eq. \(17\)](#) can be recovered by setting $\hat{\mathbf{x}} = \mathbf{B}^{-1/2}\hat{\mathbf{w}}$.

3 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^\top . 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 $s_i = \mathbf{r}_1^\top \mathbf{x}_i$, which is the coefficient of \mathbf{x}_i projected onto \mathbf{r}_1 . The variance in the direction of \mathbf{r}_1 is

$$\sum_{i=1}^n (s_i - \bar{s})^2 = \sum_{i=1}^n (\mathbf{r}_1^\top (\mathbf{x}_i - \bar{\mathbf{x}}))^2 \quad (19)$$

where $\bar{\mathbf{x}}$ is the mean over all observations. We want to find \mathbf{r}_1 such that the variance in [Eq. \(19\)](#) becomes maximal. Note that

$$\begin{aligned} \sum_{i=1}^n (\mathbf{r}_1^\top (\mathbf{x}_i - \bar{\mathbf{x}}))^2 &= \sum_{i=1}^n \mathbf{r}_1^\top (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{r}_1 \\ &= \mathbf{r}_1^\top \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{r}_1 \\ &= (n-1) \mathbf{r}_1^\top \hat{\Sigma} \mathbf{r}_1 \end{aligned} \quad (20)$$

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. \(20\)](#) is equivalent to solving the basic Rayleigh quotient maximisation problem in [Section 2](#). 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. \(19\)](#)) 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

$$\mathbf{P}^2 = \mathbf{P} = \mathbf{P}^\top. \quad (21)$$

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

$$\mathbf{P} = \sum_{i=1}^{m-1} \mathbf{r}_i \mathbf{r}_i^\top \quad (22)$$

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

$$\mathbf{X}_{m-1} = \mathbf{X} \left(\mathbf{I}_p - \sum_{i=1}^{m-1} \mathbf{r}_i \mathbf{r}_i^\top \right) \quad (23)$$

and

$$\hat{\boldsymbol{\Sigma}}_{m-1} = \frac{\mathbf{X}_{m-1}^\top \mathbf{X}_{m-1}}{n-1} = \left(\mathbf{I}_p - \sum_{i=1}^{m-1} \mathbf{r}_i \mathbf{r}_i^\top \right)^\top \hat{\boldsymbol{\Sigma}} \left(\mathbf{I}_p - \sum_{i=1}^{m-1} \mathbf{r}_i \mathbf{r}_i^\top \right). \quad (24)$$

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

$$\max_{\mathbf{r}} \mathbf{r}^\top \hat{\boldsymbol{\Sigma}}_{m-1} \mathbf{r} \quad \text{subject to} \quad \mathbf{r}^\top \mathbf{r} = 1. \quad (25)$$

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

$$\begin{aligned} \mathbf{r}_m^\top \mathbf{r}_i &= \frac{1}{\lambda_m} \mathbf{r}_m^\top \hat{\boldsymbol{\Sigma}}_{m-1} \mathbf{r}_i \\ &= \frac{1}{\lambda_m} \mathbf{r}_m^\top \left(\mathbf{I}_p - \sum_{j=1}^{m-1} \mathbf{r}_j \mathbf{r}_j^\top \right)^\top \hat{\boldsymbol{\Sigma}} \underbrace{\left(\mathbf{I}_p - \sum_{j=1}^{m-1} \mathbf{r}_j \mathbf{r}_j^\top \right) \mathbf{r}_i}_{=\mathbf{0}} = 0 \end{aligned} \quad (26)$$

as well as

$$\begin{aligned} \lambda_m \mathbf{r}_m &= \hat{\boldsymbol{\Sigma}}_{m-1} \mathbf{r}_m \\ &= \left(\mathbf{I}_p - \sum_{j=1}^{m-1} \mathbf{r}_j \mathbf{r}_j^\top \right)^\top \hat{\boldsymbol{\Sigma}} \left(\mathbf{I}_p - \sum_{j=1}^{m-1} \mathbf{r}_j \mathbf{r}_j^\top \right) \mathbf{r}_m \\ &= \left(\mathbf{I}_p - \sum_{j=1}^{m-1} \mathbf{r}_j \mathbf{r}_j^\top \right) \hat{\boldsymbol{\Sigma}} \mathbf{r}_m \\ &= \hat{\boldsymbol{\Sigma}} \mathbf{r}_m - \sum_{j=1}^{m-1} \mathbf{r}_j \underbrace{\mathbf{r}_j^\top \hat{\boldsymbol{\Sigma}} \mathbf{r}_m}_{=\lambda_j \mathbf{r}_j^\top \mathbf{r}_m = \mathbf{0}} \\ &= \hat{\boldsymbol{\Sigma}} \mathbf{r}_m \end{aligned} \quad (27)$$

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

$$\hat{\Sigma} = \frac{1}{n-1} \mathbf{X}^\top \mathbf{X} \quad (28)$$

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

$$\hat{\Sigma} \mathbf{r}_i = \lambda_i \mathbf{r}_i, \|\mathbf{r}_i\| = 1, i = 1, \dots, p \quad \text{and} \quad \mathbf{r}_i^\top \mathbf{r}_j = 0, i \neq j \quad (29)$$

as well as

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \quad (30)$$

4. Set

$$\mathbf{R} = (\mathbf{r}_1, \dots, \mathbf{r}_p) \in \mathbb{R}^{p \times p} \quad \text{and} \quad \mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_p) \quad (31)$$

so that

$$\hat{\Sigma} = \mathbf{R} \mathbf{D} \mathbf{R}^\top \quad (32)$$

5. The vectors \mathbf{r}_i are the principal component directions, the projections $\mathbf{r}_i^\top \mathbf{x}$ 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

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^p \mu_i. \quad (33)$$

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

$$\text{tr}(\hat{\Sigma}) = \sum_{i=1}^p s^2(\mathbf{X}_{\cdot i}) = \sum_{i=1}^p \lambda_i, \quad (34)$$

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 1 ⁶. 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

⁶ $\sum_{i=1}^p \lambda_i / p = 1$

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.

4 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

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top \quad (35)$$

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

$$\mathbf{U}^\top \mathbf{U} = \mathbf{I}_p \quad \text{and} \quad \mathbf{V} \mathbf{V}^\top = \mathbf{V}^\top \mathbf{V} = \mathbf{I}_p \quad (36)$$

Note that if $n > p$ it cannot hold that $\mathbf{U} \mathbf{U}^\top = \mathbf{I}_n$ ⁷. 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. (36) can now be used to derive the following equations

$$\begin{aligned} \mathbf{X} \mathbf{X}^\top \mathbf{U} &= \mathbf{U} \mathbf{D} \mathbf{V}^\top \mathbf{V} \mathbf{D} \mathbf{U}^\top \mathbf{U} = \mathbf{U} \mathbf{D}^2 \\ \mathbf{X}^\top \mathbf{X} \mathbf{V} &= \mathbf{V} \mathbf{D} \mathbf{U}^\top \mathbf{U} \mathbf{D} \mathbf{V}^\top \mathbf{V} = \mathbf{V} \mathbf{D}^2 \end{aligned} \quad (37)$$

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

$$\mathbf{X} \mathbf{X}^\top \mathbf{u}_i = d_i^2 \mathbf{u}_i \quad \text{and} \quad \mathbf{X}^\top \mathbf{X} \mathbf{v}_i = d_i^2 \mathbf{v}_i, \quad i = 1, \dots, p \quad (38)$$

Since $n \geq p$, the matrix $\mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{p \times p}$ is as large as or smaller than $\mathbf{X} \mathbf{X}^\top \in \mathbb{R}^{n \times n}$. It is therefore more computationally effective to calculate \mathbf{V} by solving the p eigenvalue problems in Eq. (38) 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.

$$\mathbf{U} = \mathbf{X} \mathbf{V} \mathbf{D}^{-1}. \quad (39)$$

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 $\mathbf{X}^\top \in \mathbb{R}^{p \times n}$. 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

$$\mathbf{X}^\top = \mathbf{V} \mathbf{D} \mathbf{U}^\top. \quad (40)$$

By transposing again we get

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top. \quad (41)$$

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

⁷A 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.

4.1 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

$$\hat{\Sigma} = \frac{\mathbf{X}^\top \mathbf{X}}{n-1}. \quad (42)$$

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

$$\hat{\Sigma} = \frac{1}{n-1} \mathbf{V} \mathbf{D} \mathbf{U}^\top \mathbf{U} \mathbf{D} \mathbf{V}^\top = \mathbf{V} \left(\frac{1}{n-1} \mathbf{D}^2 \right) \mathbf{V}^\top. \quad (43)$$

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

$$\hat{\Sigma} \mathbf{V} = \mathbf{V} \left(\frac{1}{n-1} \mathbf{D}^2 \right) \quad (44)$$

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

4.2 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 $\hat{S} \subset \mathbb{R}^p$, such that

$$\hat{S} = \arg \min_S \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{P}_S \mathbf{x}_i\|_2^2 \quad (45)$$

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

$$\begin{aligned} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{P}_S \mathbf{x}_i\|_2^2 &= \sum_{i=1}^n (\mathbf{x}_i - \mathbf{P}_S \mathbf{x}_i)^\top (\mathbf{x}_i - \mathbf{P}_S \mathbf{x}_i) \\ &= \sum_{i=1}^n [\mathbf{x}_i^\top \mathbf{x}_i - 2\mathbf{x}_i^\top \mathbf{P}_S \mathbf{x}_i + \mathbf{x}_i^\top \mathbf{P}_S^\top \mathbf{P}_S \mathbf{x}_i] \\ &= \sum_{i=1}^n [\mathbf{x}_i^\top \mathbf{x}_i - \mathbf{x}_i^\top \mathbf{P}_S^\top \mathbf{P}_S \mathbf{x}_i] \\ &= \sum_{i=1}^n [\|\mathbf{x}_i\|_2^2 - \|\mathbf{P}_S \mathbf{x}_i\|_2^2] \end{aligned} \quad (46)$$

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 vectors⁸. 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 $\mathbf{P} = \mathbf{r}_1\mathbf{r}_1^\top$ and thus we are trying to maximize

$$\begin{aligned} \sum_{i=1}^n \|\mathbf{r}_1\mathbf{r}_1^\top \mathbf{x}_i\|_2^2 &= \sum_{i=1}^n (\mathbf{r}_1^\top \mathbf{x}_i)^2 \\ &= \mathbf{r}_1^\top \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right) \mathbf{r}_1 \\ &= (n-1) \mathbf{r}_1^\top \widehat{\Sigma} \mathbf{r}_1 \end{aligned} \quad (47)$$

The term we are maximizing is the Rayleigh quotient and therefore the optimal subspace is spanned by \mathbf{r}_1 , the eigenvector of $\widehat{\Sigma}$ 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} .

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 $< n$.

4.3 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

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top = \sum_{i=1}^{\min(n,p)} d_i \mathbf{u}_i \mathbf{v}_i^\top. \quad (48)$$

Each matrix $\mathbf{u}_i \mathbf{v}_i^\top$ is of rank 1 and these matrices are scaled by the singular values d_i . As we have seen above in [Section 4.2](#), 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

$$\mathbf{P}_q = \sum_{i=1}^q \mathbf{v}_i \mathbf{v}_i^\top. \quad (49)$$

⁸Since $\sum_{i=1}^n \|\mathbf{x}_i\|_2^2$ is a constant given a dataset. The only object we can control is the subspace we project into and therefore we can only change $\sum_{i=1}^n \|\mathbf{P}_S \mathbf{x}_i\|_2^2$.

Projecting \mathbf{X} on this optimal q -dimensional subspace using the projection in Eq. (49) and the representation in Eq. (48) leads to

$$\mathbf{X}_q = \mathbf{X} \mathbf{P}_q = \left(\sum_{i=1}^{\min(n,p)} d_i \mathbf{u}_i \mathbf{v}_i^\top \right) \left(\sum_{j=1}^q \mathbf{v}_j \mathbf{v}_j^\top \right) = \sum_{i=1}^q d_i \mathbf{u}_i \mathbf{v}_i^\top. \quad (50)$$

So projecting the data into this optimal q -dimensional subspace simply means to only keep components 1 to q . The approximation error in Frobenius norm⁹ is

$$\begin{aligned} \|\mathbf{X} - \mathbf{X}_q\|_F^2 &= \left\| \sum_{i=q+1}^{\min(n,p)} d_i \mathbf{u}_i \mathbf{v}_i^\top \right\|_F^2 \\ &= \text{tr} \left[\left(\sum_{i=q+1}^{\min(n,p)} d_i \mathbf{u}_i \mathbf{v}_i^\top \right)^\top \left(\sum_{j=q+1}^{\min(n,p)} d_j \mathbf{u}_j \mathbf{v}_j^\top \right) \right] \\ &= \text{tr} \left[\sum_{i,j=q+1}^{\min(n,p)} d_i d_j \mathbf{v}_i \mathbf{u}_i^\top \mathbf{u}_j \mathbf{v}_j^\top \right] \\ &= \text{tr} \left[\sum_{i,j=q+1}^{\min(n,p)} d_i d_j \mathbf{v}_i \mathbf{v}_j^\top \right] \\ &= \sum_{k=1}^{\min(n,p)} \sum_{i,j=q+1}^{\min(n,p)} d_i d_j v_{ik} v_{jk} \\ &= \sum_{i,j=q+1}^{\min(n,p)} d_i d_j \sum_{k=1}^{\min(n,p)} v_{ik} v_{jk} \\ &= \sum_{i,j=q+1}^{\min(n,p)} d_i d_j \mathbf{1}(i=j) = \sum_{i=q+1}^{\min(n,p)} d_i^2 \end{aligned} \quad (51)$$

4.4 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

$$\mathbf{y} = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon \quad (52)$$

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

$$\varepsilon \sim \text{Normal}(0, \sigma^2) \quad (53)$$

⁹This is a matrix norm defined by $\|\mathbf{X}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^p x_{ij}^2$

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

$$\min_{\beta \in \mathbb{R}^{p+1}} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \quad (54)$$

The solution to the least squares problem is given by

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad (55)$$

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

$$\begin{aligned} \hat{\beta} &= (\mathbf{V}\mathbf{D}\mathbf{U}^\top\mathbf{U}\mathbf{D}\mathbf{V}^\top)^{-1} \mathbf{V}\mathbf{D}\mathbf{U}^\top \mathbf{y} \\ &= (\mathbf{V}\mathbf{D}^2\mathbf{V}^\top)^{-1} \mathbf{V}\mathbf{D}\mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V}\mathbf{D}^{-2}\mathbf{V}^\top \mathbf{V}\mathbf{D}\mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^\top \mathbf{y} \end{aligned} \quad (56)$$

The expression

$$\mathbf{X}^+ = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^\top \quad (57)$$

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

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{U}\mathbf{D}\mathbf{V}^\top \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^\top \mathbf{y} = \mathbf{U}\mathbf{U}^\top \mathbf{y}. \quad (58)$$

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

For ridge regression the problem to be solved is

$$\min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2 \quad (59)$$

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

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}. \quad (60)$$

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

$$\hat{\beta} = \mathbf{V} (\mathbf{D}^2 + \lambda \mathbf{I}_p)^{-1} \mathbf{D}\mathbf{U}^\top \mathbf{y} = \sum_{i=1}^p \frac{d_i}{d_i^2 + \lambda} \mathbf{v}_i \mathbf{u}_i^\top \mathbf{y}. \quad (61)$$

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.

5 Factor Analysis (FA)

Factor Analysis (FA) is a simple model that helps in the exploratory data analysis of measurements with many variables. The idea is that there are underlying unobserved factors that can explain the other variables. A major driving force behind the development of FA was to understand and make measurable abstract properties such as intelligence in psychology.

The traditional model description of FA is a linear model

$$\mathbf{x}_i = \boldsymbol{\mu} + \mathbf{L}\mathbf{z}_i + \mathbf{e}_i, \quad i = 1, \dots, n, \quad (62)$$

where

- $\mathbf{x}_i \in \mathbb{R}^p$ is the observed data,
- $\boldsymbol{\mu} \in \mathbb{R}^p$ are the mean values of each variable,
- $\mathbf{L} \in \mathbb{R}^{p \times k}$ is called the factor loading matrix describing how much each of the k factors contributes to the description of the p variables,
- $\mathbf{z}_i \in \mathbb{R}^k$ are the latent scores, describing the extent to which each p -dimensional factor describes observation \mathbf{x}_i ,
- and \mathbf{e}_i are residual errors.

Note that the mean values and the factor loading matrix are constant across observations. Denote $\mathbf{Z} = (\mathbf{z}_1 \dots \mathbf{z}_n)$ and $\mathbf{E} = (\mathbf{e}_1 \dots \mathbf{e}_n)$. Typical assumptions are

- \mathbf{Z} and \mathbf{E} are independent,
- $\mathbb{E}[\mathbf{e}_i] = \mathbf{0}$ for all i and $\text{Cov}[\mathbf{E}] = \boldsymbol{\Psi} \in \mathbb{R}^{p \times p}$ is diagonal,
- and $\mathbb{E}[\mathbf{z}_i] = \mathbf{0}$ for all i and $\text{Cov}[\mathbf{Z}] = \mathbf{I}_k$.

The first assumption ensures that all structural information is extracted and that \mathbf{E} represents the residual error, unexplained by the model. The second assumption, enforces the interpretation of \mathbf{E} as the residual error, especially by ensuring that the covariance is diagonal. This means that any covariance information is supposed to be explained by the factors. The third assumption is mostly for notational convenience. If the mean of the latent scores were not zero, $\boldsymbol{\mu}_0$ say, then $\mathbf{L}\boldsymbol{\mu}_0$ could be absorbed into $\boldsymbol{\mu}$. Similarly for the covariance matrix of \mathbf{Z} . If $\text{Cov}(\mathbf{Z}) = \mathbf{L}_0$ then setting $\tilde{\mathbf{L}} = \mathbf{L}\mathbf{L}_0^{-1/2}$ would lead to an equivalent model with unit covariance matrix (see also [Murphy \(2012, Section 12.1.1\)](#)).

Denoting $\mathbf{X} = (\mathbf{x}_1 \dots \mathbf{x}_n)$, the assumptions above ensure that

$$\mathbb{E}(\mathbf{x}_i) = \boldsymbol{\mu} \quad (63)$$

$$\text{Cov}(\mathbf{X}) = \text{Cov}(\mathbf{L}\mathbf{Z} + \mathbf{E}) = \mathbf{L}\mathbf{L}^\top + \boldsymbol{\Psi} \quad (64)$$

where in the last equation we used the independence of \mathbf{Z} and \mathbf{E} as well as the assumptions on the covariances. This shows that factor analysis provides a low-rank parametrisation¹⁰ of the covariance of the variables in \mathbf{X} . Instead of p^2 parameters only $pk + p$ are needed. Depending on the relative sizes of p and k this can lead to a substantial reduction in the number of necessary parameters.

¹⁰The rank of \mathbf{L} and therefore $\mathbf{L}\mathbf{L}^\top$ is at most k , if $k < p$.

Note that from a probabilistic perspective, the model and assumptions above are equivalent to the hierarchical model

$$\begin{aligned} \mathbf{z}_i &\sim \text{Normal}(\mathbf{0}, \mathbf{I}_k) \\ \mathbf{x}_i | \mathbf{z}_i &\sim \text{Normal}(\mathbf{L}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}) \end{aligned} \quad (65)$$

Computing the marginal distribution of \mathbf{x}_i reveals the low-rank covariance parametrisation above, i.e.

$$\mathbf{x}_i \sim \text{Normal}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^\top + \boldsymbol{\Psi}). \quad (66)$$

As mentioned in the beginning, a major driving force in the advent of FA was to find interpretable latent factors that can explain the observed variables. The posterior of $\mathbf{z}_i | \mathbf{x}_i$ is therefore of interest. Using Bayes rule together with the model in Eq. (65) we can compute¹¹

$$\mathbf{z}_i | \mathbf{x}_i \sim \text{Normal}(\mathbf{m}_i, \boldsymbol{\Sigma}) \quad (67)$$

where $\boldsymbol{\Sigma} = (\mathbf{I}_k + \mathbf{L}^\top \boldsymbol{\Psi}^{-1} \mathbf{L})^{-1}$ and $\mathbf{m}_i = \boldsymbol{\Sigma} \mathbf{L}^\top \boldsymbol{\Psi}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})$. Note that a priori, neither \mathbf{L} , $\boldsymbol{\Psi}$, nor $\boldsymbol{\mu}$ are known.

6 Probabilistic Principal Component Analysis (PPCA)

Probabilistic Principal Component Analysis (PPCA) is a variant of FA and is related to standard PCA described above. Setting $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}_p$ and choosing k as the number of principal components to retain, the model becomes very similar to PCA with added homogeneous independent Gaussian noise. Note that the predictive distribution for the latent factors can be written as

$$\mathbf{z}_i | \mathbf{x}_i \sim \text{Normal}(\mathbf{M}^{-1} \mathbf{L}^\top (\mathbf{x}_i - \boldsymbol{\mu}), \sigma^2 \mathbf{M}^{-1}) \quad (68)$$

where $\mathbf{M} = \sigma^2 \mathbf{I}_k + \mathbf{L}^\top \mathbf{L}$. It is easily seen that when $\sigma \rightarrow 0$ then the PPCA model becomes the standard PCA model, since the covariance vanishes¹² and the mean becomes $(\mathbf{L}^\top \mathbf{L})^{-1} \mathbf{L}^\top (\mathbf{x}_i - \boldsymbol{\mu})$ which is the orthogonal projection onto the column space of \mathbf{L} .

It can be shown that the maximum likelihood estimators of \mathbf{L} and σ^2 recover the usual interpretation of PCA, in the sense that if $\mathbf{L} \in \mathbb{R}^{n \times k}$ then \mathbf{L} will be estimated to a

¹¹Using Bayes rule and calculating on the log-scale

$$\begin{aligned} \log p(\mathbf{z}_i | \mathbf{x}_i) &= \log p(\mathbf{x}_i | \mathbf{z}_i) + \log p(\mathbf{z}_i) + \text{const} \\ &= \log \text{Normal}(\mathbf{x}_i; \mathbf{L}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}) + \log \text{Normal}(\mathbf{z}_i; \mathbf{0}, \mathbf{I}_k) + \text{const} \\ &= -\mathbf{z}_i \mathbf{L}^\top \boldsymbol{\Psi}^{-1} \mathbf{L} \mathbf{z}_i + 2(\mathbf{x}_i - \boldsymbol{\mu})^\top \boldsymbol{\Psi}^{-1} \mathbf{L} \mathbf{z}_i - \mathbf{z}_i^\top \mathbf{z}_i + \text{const} \\ &= -\mathbf{z}_i (\mathbf{I}_k + \mathbf{L}^\top \boldsymbol{\Psi}^{-1} \mathbf{L}) \mathbf{z}_i \\ &\quad + 2 \left((\mathbf{I}_k + \mathbf{L}^\top \boldsymbol{\Psi}^{-1} \mathbf{L})^{-1} \mathbf{L}^\top \boldsymbol{\Psi}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right)^\top (\mathbf{I}_k + \mathbf{L}^\top \boldsymbol{\Psi}^{-1} \mathbf{L}) \mathbf{z}_i + \text{const} \end{aligned}$$

where the symmetry of $\boldsymbol{\Psi}$ was used.

¹²i.e. the model becomes degenerate/singular

slightly penalised version of the first k principal components and σ^2 is the average of the remaining singular values.

There are multiple advantages to the probabilistic formulation of PCA. One advantage of the probabilistic formulation of PCA is that it allows for data driven selection of the number of latent components to retain. One possible solution for this is Automatic Relevance Determination (ARD). Another advantage is that it is easier to adapt PCA to other data types than continuous data as well as the incorporation of missing data (Roweis, 1998; Tipping and Bishop, 1999).

7 Canonical Correlation Analysis (CCA)

The previous methods focused only on one dataset and projections into lower dimensions that emphasize certain aspects (e.g. variance or correlations) of the data. When two datasets are analysed jointly, one question that can be asked is how to project each dataset into a lower-dimensional space such that the most similar directions in each dataset are uncovered.

The traditional view/interpretation of CCA is as follows: Given two samples of (possibly different) features, each sample vector of dimension p and q , respectively, find all pairs of directions in the respective feature space (i.e. linear combinations of the feature vectors) that are maximally correlated and orthogonal to all previous found pairs.

This can be formalised as follows. Assume that there are data matrices $\mathbf{X} \in \mathbb{R}^{n \times p}$ and $\mathbf{Y} \in \mathbb{R}^{n \times q}$. The task of CCA in its first step is to find $\mathbf{a} \in \mathbb{R}^p$ and $\mathbf{b} \in \mathbb{R}^q$ such that the sample correlation between vectors $\mathbf{X}\mathbf{a}$ and $\mathbf{Y}\mathbf{b}$ is maximal. Let $\hat{\Sigma}_{XY}$ be the cross-covariance matrix between \mathbf{X} and \mathbf{Y} , and $\hat{\Sigma}_{XX}$ as well as $\hat{\Sigma}_{YY}$ the respective covariance matrices. A necessary assumption is that the covariance matrices of \mathbf{X} and \mathbf{Y} are of full rank, i.e. are positive definite. The objective function is then

$$r(\mathbf{a}, \mathbf{b}) = s_{\mathbf{X}\mathbf{a}, \mathbf{Y}\mathbf{b}} = \frac{\mathbf{a}^\top \hat{\Sigma}_{XY} \mathbf{b}}{\left(\mathbf{a}^\top \hat{\Sigma}_{XX} \mathbf{a}\right)^{1/2} \left(\mathbf{b}^\top \hat{\Sigma}_{YY} \mathbf{b}\right)^{1/2}}. \quad (69)$$

Note that $r(\mathbf{a}, \mathbf{b})$ is invariant under scaling of either argument and therefore we can impose the restrictions

$$\mathbf{a}^\top \hat{\Sigma}_{XX} \mathbf{a} = 1 \quad \text{and} \quad \mathbf{b}^\top \hat{\Sigma}_{YY} \mathbf{b} = 1. \quad (70)$$

Using the variable substitutions $\mathbf{c} \leftarrow \hat{\Sigma}_{XX}^{1/2} \mathbf{a}$ and $\mathbf{d} \leftarrow \hat{\Sigma}_{YY}^{1/2} \mathbf{b}$ the objective function is transformed to

$$r(\mathbf{a}, \mathbf{b}) = \mathbf{c}^\top \hat{\Sigma}_{XX}^{-1/2} \hat{\Sigma}_{XY} \hat{\Sigma}_{YY}^{-1/2} \mathbf{d}. \quad (71)$$

Denote $\mathbf{K} = \hat{\Sigma}_{XX}^{-1/2} \hat{\Sigma}_{XY} \hat{\Sigma}_{YY}^{-1/2}$ and the optimization problem to be solved, formulated as a minimization problem, becomes

$$\inf_{\mathbf{c}, \mathbf{d}} -\mathbf{c}^\top \mathbf{K} \mathbf{d} \quad \text{such that} \quad \mathbf{c}^\top \mathbf{c} = 1, \mathbf{d}^\top \mathbf{d} = 1. \quad (72)$$

where \mathbf{a} and \mathbf{b} can be recovered by re-substitution. The Lagrangian for the corresponding minimization problem is

$$L(\mathbf{c}, \mathbf{d}, \lambda) = -\mathbf{c}^\top \mathbf{K} \mathbf{d} + \lambda_1 (\mathbf{c}^\top \mathbf{c} - 1) + \lambda_2 (\mathbf{d}^\top \mathbf{d} - 1). \quad (73)$$

The dual function of the Lagrangian is

$$g(\lambda) = \inf_{\mathbf{c}, \mathbf{d}} -\mathbf{c}^\top \mathbf{K} \mathbf{d} + \lambda_1 (\mathbf{c}^\top \mathbf{c} - 1) + \lambda_2 (\mathbf{d}^\top \mathbf{d} - 1) \quad (74)$$

and can be rewritten as either

$$g(\lambda) = \inf_{\mathbf{c}, \mathbf{d}} \frac{1}{4\lambda_1} \left\{ \|2\lambda_1 \mathbf{c} - \mathbf{K} \mathbf{d}\|_2^2 + \mathbf{d}^\top (-\mathbf{K}^\top \mathbf{K} + 4\lambda_1 \lambda_2 \mathbf{I}_q) \mathbf{d} \right\} - \lambda_1 - \lambda_2, \quad (75)$$

or

$$g(\lambda) = \inf_{\mathbf{c}, \mathbf{d}} \frac{1}{4\lambda_2} \left\{ \|2\lambda_2 \mathbf{d} - \mathbf{K}^\top \mathbf{c}\|_2^2 + \mathbf{c}^\top (-\mathbf{K} \mathbf{K}^\top + 4\lambda_1 \lambda_2 \mathbf{I}_p) \mathbf{c} \right\} - \lambda_1 - \lambda_2. \quad (76)$$

The reformulation in Eq. (75) shows directly that if $\lambda_1 < 0$ then the infimum is $-\infty$ by choosing $\mathbf{d} = \mathbf{0}$ and letting $\|\mathbf{c}\|_2 \rightarrow \infty$. An analogous result for $\lambda_2 < 0$ follows from Eq. (76). It follows from Eq. (74) that the infimum is $-\infty$ if $\lambda_1 = 0$ by choosing $\mathbf{c} = \tau \mathbf{K} \mathbf{d}$ and letting $\tau \rightarrow \infty$ for any \mathbf{d} such that $\mathbf{K} \mathbf{d} \neq \mathbf{0}$ (exists as long as $\mathbf{K} \neq \mathbf{0}$). An analogous result follows for $\lambda_2 = 0$.

In the case $\lambda_1 > 0$ and $\lambda_2 > 0$, it is enough to focus on Eq. (75) if $p > q$ or Eq. (76) if $q > p$. Assume $p > q$ here¹³. The arguments are similar to the analysis of the dual function of the basic Rayleigh coefficient.

If the matrix $(-\mathbf{K}^\top \mathbf{K} + 4\lambda_1 \lambda_2 \mathbf{I}_q)$ is not psd there exists a negative eigenvalue μ with corresponding eigenvector \mathbf{u} . Setting $\mathbf{d} = \tau \mathbf{u}$ and $\mathbf{c} = \frac{1}{2\lambda_1} \mathbf{K} \mathbf{d}$ leads to

$$\|2\lambda_1 \mathbf{c} - \mathbf{K} \mathbf{d}\|_2^2 + \mathbf{d}^\top (-\mathbf{K}^\top \mathbf{K} + 4\lambda_1 \lambda_2 \mathbf{I}_q) \mathbf{d} = \tau^2 \mu \|\mathbf{u}\|_2^2 \rightarrow -\infty \quad \text{for } \tau \rightarrow \infty. \quad (77)$$

For the infimum to be greater than $-\infty$ it is therefore necessary that $(-\mathbf{K}^\top \mathbf{K} + 4\lambda_1 \lambda_2 \mathbf{I}_q)$ is psd. Let the SVD of \mathbf{K} be given by $\mathbf{K} = \mathbf{U} \mathbf{\Delta} \mathbf{V}^\top$ assuming that the singular values $\delta_1 \geq \delta_2 \geq \dots \geq \delta_q \geq 0$ in $\mathbf{\Delta}$ are sorted from largest to smallest. Denote the columns of \mathbf{V} by \mathbf{v}_i and the columns of \mathbf{U} by \mathbf{u}_i . Since $p > q$ it holds that $\mathbf{V} \mathbf{V}^\top = \mathbf{I}_q$ which implies

$$\begin{aligned} \mathbf{d}^\top (-\mathbf{K}^\top \mathbf{K} + 4\lambda_1 \lambda_2 \mathbf{I}_q) \mathbf{d} &= \mathbf{d}^\top \mathbf{V} (-\mathbf{\Delta}^2 + 4\lambda_1 \lambda_2 \mathbf{I}_q) \mathbf{V}^\top \mathbf{d} \\ &= \sum_{i=1}^q (4\lambda_1 \lambda_2 - \delta_i^2) (\mathbf{v}_i^\top \mathbf{d})^2 \end{aligned} \quad (78)$$

for all $\mathbf{d} \in \mathbb{R}^q$. The matrix $(-\mathbf{K}^\top \mathbf{K} + 4\lambda_1 \lambda_2 \mathbf{I}_q)$ is therefore psd if $4\lambda_1 \lambda_2 \geq \delta_1^2$. The norm and quadratic terms in Eq. (75) are therefore greater or equal zero and both become zero for $\mathbf{c} = \mathbf{0} = \mathbf{d}$. The dual function is therefore

$$g(\lambda) = \begin{cases} -\lambda_1 - \lambda_2 & \text{if } 4\lambda_1 \lambda_2 \geq \delta_1^2, \lambda_1 > 0, \lambda_2 > 0, \\ -\infty & \text{otherwise} \end{cases} \quad (79)$$

¹³The arguments for $q > p$ are symmetric using Eq. (76) instead.

and implies the dual problem

$$\sup_{\lambda} -\lambda_1 - \lambda_2 \quad \text{such that} \quad 4\lambda_1\lambda_2 \geq \delta_1^2, \lambda_1 > 0, \lambda_2 > 0. \quad (80)$$

It can be easily seen¹⁴ that the supremum is achieved for $\hat{\lambda}_1 = \delta_1/2 = \hat{\lambda}_2$ which leads to a dual objective function value of $d^* = -\delta_1$.

Let $\mathbf{c}_1 = \mathbf{u}_1$ and $\mathbf{d}_1 = \mathbf{v}_1$. Since the primal objective value p^* in Eq. (72) is bounded from below by d^* , it holds that

$$d^* = -\delta_1 \leq p^* \leq -\mathbf{u}_1^\top \mathbf{K} \mathbf{v}_1 = -\delta_1. \quad (81)$$

This shows that $p^* = -\delta_1$ and that the optimisation problem in Eq. (72) is solved by \mathbf{c}_1 and \mathbf{d}_1 . By resubstitution it follows that $\mathbf{a}_1 = \hat{\Sigma}_{XX}^{-1/2} \mathbf{u}_1$ and $\mathbf{b}_1 = \hat{\Sigma}_{YY}^{-1/2} \mathbf{v}_1$ solve the original optimisation problem of maximizing $r(\mathbf{a}_1, \mathbf{b}_1)$ in Eq. (69).

By restricting the search space for \mathbf{c} to vectors that are orthogonal to $\mathbf{c}_1, \dots, \mathbf{c}_{i-1}$ and for \mathbf{d} to vectors that are orthogonal to $\mathbf{d}_1, \dots, \mathbf{d}_{i-1}$ it follows easily—by the arguments above—that the restriction in the dual problem becomes $4\lambda_1\lambda_2 \geq \delta_i^2$ which ultimately leads to the solution $p^* = -\delta_i$ which is achieved for $\mathbf{c}_i = \mathbf{u}_i$ and $\mathbf{d}_i = \mathbf{v}_i$. By defining $\boldsymbol{\eta}_i = \mathbf{X} \mathbf{a}_i$ and $\boldsymbol{\rho}_i = \mathbf{Y} \mathbf{b}_i$, q pairs¹⁵ of canonical directions $(\boldsymbol{\eta}_1, \boldsymbol{\rho}_1), \dots, (\boldsymbol{\eta}_q, \boldsymbol{\rho}_q)$ can be found such that

$$\hat{\Sigma}_{(\boldsymbol{\eta}, \boldsymbol{\rho})(\boldsymbol{\eta}, \boldsymbol{\rho})} = \begin{pmatrix} \mathbf{I}_p & \boldsymbol{\Delta} \\ \boldsymbol{\Delta} & \mathbf{I}_q \end{pmatrix}. \quad (82)$$

CCA can also be interpreted as a probabilistic model as shown by Bach and Jordan (2006).

8 Independent Component Analysis (ICA)

References

- F. R. Bach and M. I. Jordan. A Probabilistic Interpretation of Canonical Correlation Analysis. Technical Report 688, Department of Statistics, University of California, Berkeley, 2006. (page 17)
- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004. (page 2)
- M. Falk, F. Marohn, and B. Tewes. *Foundations of Statistical Analyses and Applications with SAS*. Birkhäuser, 2002. (page 2)
- W. Härdle and L. Simar. *Applied Multivariate Statistical Analysis*. Springer-Verlag Berlin Heidelberg, 2 edition, 2007. (page 2)

¹⁴TODO Create graphic

¹⁵If $q > p$, then only p pairs can be found.

- K. P. Murphy. *Machine Learning: A Probabilistic Perspective*. Massachusetts Institute of Technology, 2012. (pages [2](#), [13](#))
- S. T. Roweis. EM Algorithms for PCA and SPCA. In M. I. Jordan, M. J. Kearns, and S. A. Solla, editors, *Advances in Neural Information Processing Systems 10*, pages 626–632. MIT Press, 1998. (page [15](#))
- M. E. Tipping and C. M. Bishop. Probabilistic principal component analysis. *J. R. Statist. Soc. B*, 61(Part 3):611–622, 1999. (page [15](#))