

Multivariate Statistical Analysis

Selected Lecture Notes

Radoslav Harman
Faculty of Mathematics, Physics and Informatics, Comenius University

May 17, 2026

Contents

1	Principal Components Analysis	2
1.1	Mathematical background	2
1.2	Theoretical principal components	4
1.3	Sample principal components	8
1.4	History and applications of principal components	10
2	Multidimensional Scaling	12
2.1	Theory for Classical Multidimensional Scaling	12
2.2	Application of Classical Multidimensional Scaling	14
3	Canonical Correlations	15
3.1	Mathematical background for Canonical Correlations	15
3.2	Theoretical Canonical Correlations	15
3.3	Sample Canonical Correlations	19
3.4	Applications of Canonical Correlations	20
4	Factor Analysis	20
4.1	Estimation of factor loadings	22
4.2	Rotation of Factors	23
4.3	Estimation of Factor Scores	24
5	Partitioning Methods of Cluster Analysis	25
5.1	Clustering using k -means and k -medoids	27
5.2	Model-based Clustering	30
5.3	Clustering using DBSCAN	35
5.4	Using clustering for anomaly detection	37
6	Classification methods in general	38
6.1	Representations of classifiers	39
6.2	Reliability of a classifier	40
6.3	Binary classification	42
6.4	Estimation of misclassification rates	44
7	Bayes classifier	46
7.1	The general Bayes classifier	46
7.2	Naive Bayes classifier	49
7.3	Linear and quadratic discriminant analyses	50
7.4	K nearest neighbours	53

8	Classification trees	54
8.1	Constructing a classification tree using recursive partitioning	54
8.2	Pruning of a classification tree	57
8.3	Bagging and classification forests	58
8.4	Boosting	59
9	Support vector machines	62
9.1	Linear support vector machines	62
9.2	Nonlinear support vector machines	65
10	Important properties of classifiers	68

1 Principal Components Analysis

1.1 Mathematical background

We assume that the reader is already familiar with fundamental notions and results of matrix algebra and multivariate probability, but we will give a brief review of some of the facts that are particularly important for the multivariate statistics methods presented in this textbook.

Recall that a $p \times p$ matrix Σ is **non-negative definite**¹, if it is symmetric and satisfies $a^T \Sigma a \geq 0$ for any vector $a \in \mathbb{R}^p$. If $\Sigma u = \lambda u$ for some $\lambda \in \mathbb{R}$ and $u \in \mathbb{R}^p$, $u \neq 0$, then u is an **eigenvector** of Σ and λ is the **eigenvalue** of Σ that corresponds to u .

Vectors u_1, \dots, u_p form an **orthonormal system** if these vectors are mutually orthogonal and each is normalized such that $\|u_i\| = 1$ for all $i = 1, \dots, p$. Matrix U of the type $p \times p$ is an **orthogonal matrix** if $UU^T = I_p$, where I_p denotes the $p \times p$ identity matrix. That is, U is orthogonal if and only if the columns of U form an orthonormal system of vectors. The linear mapping corresponding to an orthogonal matrix is a rotation or a composition of reflection and rotation.

Theorem 1.1 (Spectral decomposition of a non-negative definite matrix). *For any non-negative definite $p \times p$ matrix Σ there exists an orthonormal system u_1, \dots, u_p of eigenvectors such that*

$$\Sigma = \sum_{i=1}^p \lambda_i u_i u_i^T = U \Lambda U^T, \tag{1}$$

¹A non-negative definite matrix is sometimes called “positive semidefinite”.

where λ_i is the eigenvalue of Σ corresponding to the eigenvector u_i for all $i = 1, \dots, p$, $U = (u_1, \dots, u_p)$ is the orthogonal matrix of normalized eigenvectors and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ is the diagonal matrix with the eigenvalues on the diagonal. If $\lambda_1 > \lambda_2 > \dots > \lambda_p$, then the eigenvectors u_1, \dots, u_p are uniquely defined (up to a possible change of the sign).

A $p \times p$ matrix Σ is **positive definite**, if it is symmetric and satisfies $a^T \Sigma a > 0$ for any vector $0 \neq a \in \mathbb{R}^p$. A matrix Σ is positive definite if and only if Σ is a non-negative definite and non-singular; equivalently, all eigenvalues of Σ are strictly positive.

An **orthogonal projector** onto a k -dimensional linear space $\mathcal{A} \subseteq \mathbb{R}^p$ is the unique symmetric matrix P of the type $p \times p$ such that $Py \in \mathcal{A}$ for all $y \in \mathbb{R}^p$, $Px = x$ for all $x \in \mathcal{A}$, and $y - Py$ is orthogonal to Pz for all $y, z \in \mathbb{R}^p$, which we denote $(y - Py) \perp Pz$. If A is a $p \times k$ matrix with rank k , where $k \leq p$, then $A^T A$ is a non-singular matrix and $P = A(A^T A)^{-1} A^T$ is the orthogonal projector on the linear space $\mathcal{C}(A)$, which is the space generated by the columns of A .

For a p -dimensional random vector $\mathbf{X} = (X_1, \dots, X_p)^T$, the **variance-covariance matrix** is a $p \times p$ matrix Σ with elements $\Sigma_{ij} = \text{cov}(X_i, X_j)$, $i, j = 1, \dots, p$. The variance-covariance matrix is always non-negative definite and, typically², it is also non-singular, i.e., positive definite. Geometrically, Σ determines the “shape” of the multivariate data generated as independent samples of \mathbf{X} .

More generally, by *Cov* we will denote the matrix of all mutual covariances of the components of a pair of random vectors \mathbf{X} and \mathbf{Z} . For multivariate statistical analysis, it is very important to know how the variance-covariance matrix and the matrix *Cov* change under linear transformations of the random vector(s).

Theorem 1.2 (The effect of a linear transformation on the variance-covariance matrix). *If \mathbf{X} is a p -dimensional random vector with covariance matrix Σ and A is an $m \times p$ matrix, then the m -dimensional random vector $\mathbf{Y} = A\mathbf{X}$ has variance-covariance matrix $A\Sigma A^T$. More generally: If \mathbf{X} is a p -dimensional random vector, \mathbf{Z} is an r -dimensional random vector, A is an $m \times p$ matrix and B is an $k \times r$ matrix, then $\text{Cov}(A\mathbf{X}, B\mathbf{Z}) = A\text{Cov}(\mathbf{X}, \mathbf{Z})B^T$.*

Principal components are based on a rotation (i.e., a specific linear transformation) of an underlying random vector as detailed in the next section.

²For instance if the random vector \mathbf{X} has a distribution continuous with respect to the Lebesgue measure in \mathbb{R}^p .

1.2 Theoretical principal components

Let μ be the mean value vector and let Σ be the variance-covariance matrix of a random vector $\mathbf{X} = (X_1, \dots, X_p)^T$ which corresponds to p -dimensional measurements or observations on n objects³.

Let u_1, \dots, u_p be an orthonormal system of eigenvectors of Σ , let $\lambda_1 \geq \dots \geq \lambda_p$ be the corresponding eigenvalues and let $U = (u_1, \dots, u_p)$; cf. Theorem 1.1. Vectors u_1, \dots, u_p determine what we will call “principal variance directions” and $\lambda_1, \dots, \lambda_p$ determine the variances in the principal directions.

Principal variance directions can be illustrated on an example of a random vector with a two-dimensional normal distribution. For instance, if

$$\Sigma = \begin{pmatrix} 3.25 & 1.30 \\ 1.30 & 1.75 \end{pmatrix},$$

then $u_1 \approx (0.87; 0.50)^T$, $u_2 \approx (-0.50; 0.87)^T$, $\lambda_1 \approx 4$ and $\lambda_2 \approx 1$. In Figure 1 we see points randomly generated from $N_2((0, 0)^T, \Sigma)$. Arrows denote the directions of vectors u_1 and u_2 . The length of the arrows is proportional to $\sqrt{\lambda_1}$ and $\sqrt{\lambda_2}$.

Clearly, eigenvectors and eigenvalues of the variance-covariance matrix Σ capture important aspects of the “shape” of the distribution of \mathbf{X} . The essence of principal components analysis (PCA) is the rotation of \mathbf{X} (or a random sample), to the coordinate system determined by the eigenvectors of Σ (or by the eigenvectors of the sample variance-covariance matrix S_n , see Subsection 1.3).

An immediate consequence of Theorems 1.1 and 1.2 is:

Theorem 1.3 (De-correlation of a random vector). *Random vector $\mathbf{Y} = U^T(\mathbf{X} - \mu)$ has the zero mean value and its variance-covariance matrix is $D(\mathbf{Y}) = \text{diag}(\lambda_1, \dots, \lambda_p)$. That is, the components of \mathbf{Y} are uncorrelated, their variances are $\lambda_1 \geq \dots \geq \lambda_p$ and their standard deviations are $\sqrt{\lambda_1} \geq \dots \geq \sqrt{\lambda_p}$.*

The previous theorem gives a theoretical basis for the following definition.

³We will only consider random variables with finite mean values and variances, that is, we will assume that all random vectors have well-defined, finite variance-covariance matrices.

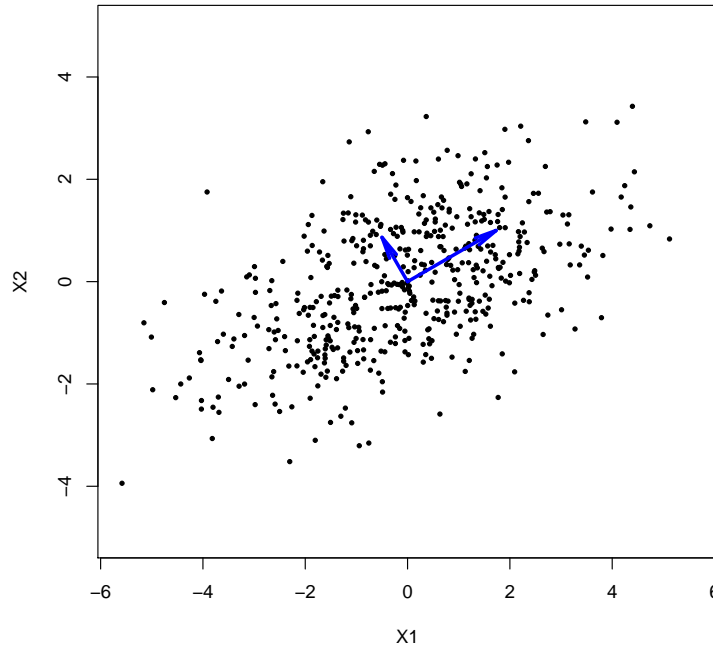


Figure 1: Principal variance directions defined by a pair of orthogonal eigenvectors of a two-dimensional normal distribution. In the figure, the lengths of the eigenvectors are proportional to the standard deviation of the corresponding principal components.

Definition 1.1 (Principal components of a random vector). *Random vector $\mathbf{Y} = (Y_1, \dots, Y_p)^T$ from Theorem 1.3 is called the **vector of (theoretical) principal components** of the random vector \mathbf{X} . For $i = 1, \dots, p$, the random variable $Y_i = u_i^T(\mathbf{X} - \mu)$ is called the **i -th principal component** of the random vector \mathbf{X} .*⁴

It is simple to show that for all $i \in \{1, \dots, p\}$ the random vector $Y_i u_i$ is the orthogonal projection of $\mathbf{X} - \mu$ onto the 1-dimensional subspace (i.e., a line) defined by all real multiples of u_i . In other words, principal components Y_1, \dots, Y_p form (random) coordinates of $\mathbf{X} - \mu$ in the coordinate system defined by the orthonormal vectors u_1, \dots, u_p .

Theorem 1.3 states that principal components Y_1, \dots, Y_p of a random

⁴Note that if some of the eigenvalues of Σ are equal, then there is an infinite number of possible choices of the corresponding eigenvectors, i.e., feasible definitions of the vectors of principal components.

vector \mathbf{X} are uncorrelated. Thus, the transformation to principal components is occasionally referred to as “de-correlation” of \mathbf{X} . Importantly, variances of principal components are in decreasing order. Note also that the sum of variances of principal components is the same as the sum of variances of variables X_1, \dots, X_p which is probably the basis for the expression that (all) principal components “explain” (all) variation in the data.

Mutual relations of the coordinates of the original random vector \mathbf{X} and the principal components of \mathbf{X} are given by the following theorem.

Theorem 1.4 (Relations of original variables and principal components). *Let $\mathbf{Y} = (Y_1, \dots, Y_p)^T$ be the vector of principal components of the random vector \mathbf{X} from Theorem 1.3. Then, for any pair $i, j \in \{1, \dots, p\}$, we have*

$$\text{cov}(X_i, Y_j) = u_{ij}\lambda_j, \quad \rho(X_i, Y_j) = u_{ij}\sqrt{\lambda_j}/\sigma_i, \quad (2)$$

where $u_{ij} = (u_j)_i$ is the i, j -th element of the matrix U (that is, the i -th coordinate of the eigenvector u_j), and $\sigma_i = \sqrt{DX_i}$.

Proof. From $X_i = e_i^T \mathbf{X}$, where e_i is the i -th standard unit vector and from $Y_j = u_j^T (\mathbf{X} - \mu)$ we obtain

$$\begin{aligned} \text{cov}(X_i, Y_j) &= \text{cov}(e_i^T \mathbf{X}, u_j^T (\mathbf{X} - \mu)) = e_i^T D(\mathbf{X}) u_j = \\ &= e_i^T \left(\sum_{k=1}^p \lambda_k u_k u_k^T \right) u_j =^* e_i^T \lambda_j u_j = \lambda_j e_i^T u_j = \lambda_j u_{ij}. \end{aligned}$$

The equality denoted by the asterisk follows from the fact that u_1, \dots, u_p are mutually orthogonal and have the unit length. \square

Exercise 1.1. *Prove the following claim. Let $b_1, b_2 \in \mathbb{R}^p$, $b_1, b_2 \neq 0$ and let \mathbf{X} be a p -dimensional random vector with positive definite covariance matrix. Then the random variables $b_1^T \mathbf{X}$ and $b_2^T \mathbf{X}$ are uncorrelated if and only if b_1 is orthogonal to b_2 in the Σ -inner product, i.e., $b_1^T \Sigma b_2 = 0$.*

The following theorem provides an important optimization/probabilistic interpretation of principal components.

Theorem 1.5 (Maximum variance justification of principal components). *Assume that all eigenvalues $\lambda_1, \dots, \lambda_p$ of Σ are distinct. The first principal component Y_1 of the p -dimensional random vector \mathbf{X} has the largest variance from all normed linear combinations of the components of \mathbf{X} . Formally:*

$$D(Y_1) = \max\{D(b^T \mathbf{X}) : b \in \mathbb{R}^p, \|b\| = 1\}$$

For $k \geq 2$, the k -th principal component Y_k of the random vector \mathbf{X} has the largest variance from all normed linear combinations of the components of \mathbf{X} that are uncorrelated with Y_1, \dots, Y_{k-1} . Formally⁵

$$D(Y_k) = \max\{D(b^T \mathbf{X}) : b \in \mathbb{R}^p, \|b\| = 1, b \perp u_1, \dots, b \perp u_{k-1}\}.$$

Proof. Let u_1, \dots, u_p be an orthonormal system of eigenvectors of the covariance matrix Σ of the p -dimensional random vector \mathbf{X} , and let $\lambda_1 > \dots > \lambda_p$ be the corresponding eigenvalues. Let $b \in \mathbb{R}^p$, $\|b\| = 1$, and let $b = \sum_{i=1}^p c_i u_i$ be the expression of b in the orthonormal basis u_1, \dots, u_p of the space \mathbb{R}^p . Since the vectors u_1, \dots, u_p are orthonormal, we obtain $u_j^T b = \sum_i c_i u_j^T u_i = c_j$ for any $j \in \{1, \dots, p\}$, and $\sum_i c_i^2 = \sum_i c_i^2 u_i^T u_i = \sum_i (c_i u_i)^T \sum_k (c_k u_k) = b^T b = 1$. Therefore

$$D(b^T \mathbf{X}) = b^T \left(\sum_{i=1}^p \lambda_i u_i u_i^T \right) b = \sum_{i=1}^p \lambda_i c_i^2 \leq^* \lambda_1 = D(Y_1). \quad (3)$$

The inequality denoted by the asterisk follows from the fact that $\sum_{i=1}^p c_i^2 = 1$, i.e., $\sum_{i=1}^p \lambda_i c_i^2$ is a weighted average of eigenvalues, which of course cannot be larger than the largest eigenvalue λ_1 . Since $Y_1 = u_1^T \mathbf{X}$, that is, Y_1 is itself a normed linear combination of the coordinates of \mathbf{X} , we obtain the first part of the theorem.

If we have $2 \leq k \leq p$ and an additional condition $b \perp u_1, \dots, b \perp u_{k-1}$, then $c_i = 0$ for all $i = 1, \dots, k-1$, which implies

$$D(b^T \mathbf{X}) = \sum_{i=1}^p \lambda_i c_i^2 = \sum_{i=k}^p \lambda_i c_i^2 \leq^* \lambda_k = D(Y_k),$$

in a way analogous to (3). □

The transformation of a random vector to principal components has several alternative geometric/optimization interpretations. For instance, let $\mathbf{X} \sim N_p(0, \Sigma)$ and by $P_{\mathcal{A}}$ denote the orthogonal projector on a linear space \mathcal{A} . Let $\mathcal{A}^* \subseteq \mathbb{R}^p$ be the k -dimensional plane that optimally fits the distribution of \mathbf{X} in the sense of least squares, i.e., \mathcal{A}^* minimizes $E(\|X - P_{\mathcal{A}^*} X\|^2)$. Then, it is possible to show that \mathcal{A}^* is spanned by the eigenvectors u_1, \dots, u_k of Σ corresponding to the k largest eigenvalues.

A measure of **proportion of variance** “explained” by the first k principal components Y_1, \dots, Y_k , $k \leq p$, or a “goodness of fit” measure of the

⁵See also Exercise 1.1.

k -dimensional plane \mathcal{A}^* from the previous paragraph, is the quantity

$$\alpha_k = \frac{\lambda_1 + \dots + \lambda_k}{\lambda_1 + \dots + \lambda_p}.$$

The quantity α_k is important for the selection of the number of principal components that capture significant amount of variability of the original p -dimensional data. It turns out that the variance-covariance matrix of multidimensional data is often such that α_k is close to 1 for values of k that are small relative to p . Geometrically, this means that the ellipsoid of dispersion is “thick” in a few orthogonal directions and “thin” in all others.

1.3 Sample principal components

In real applications, the mean value μ and the variance-covariance matrix Σ are rarely known, and the “theoretical” principal components from Definition 1.1 cannot be calculated. Usually, we only have a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n$, $n \geq p + 1$, from an otherwise unknown distribution, and the parameters μ, Σ need to be estimated by the **vector of mean values** $\bar{\mathbf{X}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$ and the **sample variance-covariance matrix**

$$S_n = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^T.$$

Definition 1.2 (Sample eigenvalues and eigenvectors). *Ordered eigenvalues of S_n will be called **sample eigenvalues** and denoted by $\hat{\lambda}_1^{(n)} > \dots > \hat{\lambda}_p^{(n)}$. The corresponding normalized eigenvectors of S_n will be called **sample eigenvectors** and denoted by $\hat{u}_1^{(n)}, \dots, \hat{u}_p^{(n)}$.*⁶

Note that (from the point of view before collecting the data) the sample eigenvalues $\hat{\lambda}_i^{(n)}$ are random variables and the sample eigenvectors $\hat{u}_i^{(n)}$ are random vectors. In general, $\hat{\lambda}_i^{(n)} \rightarrow \lambda_i$ and $\hat{u}_i^{(n)} \rightarrow u_i$ as $n \rightarrow \infty$, but the analysis of the stochastic convergence is usually very non-trivial. For the case of a normal sample, the convergence is given by the following result.⁷

⁶In theory, some of the eigenvalues of S_n could be equal, but for the case of an independent sample of size $n \geq p + 1$ from a continuous p -dimensional distribution, which is typical for applications, the probability of this event is zero. Moreover, even in the case of distinct eigenvalues $\hat{\lambda}_i^{(n)}$, the eigenvectors are not uniquely defined in the sense that if $\hat{u}_i^{(n)}$ is a normalized eigenvector, so is $-\hat{u}_i^{(n)}$. We assume that we have some consistent way of selecting which of these two eigenvectors is *the* normalized eigenvector. It would be counter-productive to be too mathematically rigorous at this point.

⁷Note, however, that in applications the normality is usually not required; the aim of PCA is almost always a reduction of dimensionality, with the aim to compress the data or to discover new knowledge, see Subsection 1.4, not hypothesis testing.

Theorem 1.6. *Assume that all eigenvalues $\lambda_1, \dots, \lambda_p$ of Σ are distinct. For any $i \in \{1, \dots, p\}$, we have the following convergence in distribution:*

$$\begin{aligned} \sqrt{n-1} (\hat{\lambda}_i^{(n)} - \lambda_i) &\rightarrow N(0, 2\lambda_i^2), \\ \sqrt{n-1} (\hat{u}_i^{(n)} - u_i) &\rightarrow N_p \left(0, \sum_{j \neq i} \lambda_i \lambda_j (\lambda_i - \lambda_j)^{-2} u_j u_j^T \right). \end{aligned}$$

Thus, if the sample size is large enough ($n \gg p$), the values of $\hat{\lambda}_i^{(n)}$ and $\hat{u}_i^{(n)}$ can usually be taken as reliable approximations of the eigenvalues and eigenvectors of Σ , and the “speed of convergence” of the estimators is approximately \sqrt{n} . Note, however, that the variances and covariances of the limiting distributions depend on the estimated parameters. Moreover, the variance covariance matrix of $\hat{u}_i^{(n)}$ ’s tends to be large, if some of the eigenvalues are very similar.

For actual realizations $\mathbf{x}_1, \dots, \mathbf{x}_n$ of the random vectors $\mathbf{X}_1, \dots, \mathbf{X}_n$ ⁸, it is often useful to compute the vectors $\mathbf{y}_i = (\hat{u}_1^{(n)}, \dots, \hat{u}_p^{(n)})^T (\mathbf{x}_i - \bar{\mathbf{x}})$, $i = 1, \dots, n$, which are called **principal components scores**. These vectors are estimates of the coordinates of \mathbf{x}_i , $i = 1, \dots, n$ in the orthonormal basis determined by the principal component directions.

In other words, assume that \mathcal{A}^* is the k -dimensional plane that best fits the data in the sense of least squares in the p -dimensional space⁹ and let $\mathbf{z}_1, \dots, \mathbf{z}_n$ be the orthogonal projections of the feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ onto \mathcal{A}^* . Then, the first k coordinates of the scores $\mathbf{y}_1, \dots, \mathbf{y}_n$ correspond to the coordinates of $\mathbf{z}_1, \dots, \mathbf{z}_n$ in the hyperplane \mathcal{A}^* . They can be used to represent the n objects in the k -dimensional space, usually in plane ($k = 2$).

The proportion α_k , $k \in \{1, \dots, p\}$, of the variance explained by the first k principal components can be estimated by

$$\hat{\alpha}_k^{(n)} = \frac{\hat{\lambda}_1^{(n)} + \dots + \hat{\lambda}_k^{(n)}}{\hat{\lambda}_1^{(n)} + \dots + \hat{\lambda}_p^{(n)}}.$$

For the random variables $\hat{\alpha}_k^{(n)}$, a theorem similar to 1.6 can be proved for the case of normality, stating that the random variables $\hat{\alpha}_k^{(n)}$ converge to the values α_k with a speed proportional to \sqrt{n} .

⁸The realizations $\mathbf{x}_1, \dots, \mathbf{x}_n$ are sometimes called feature vectors of objects.

⁹Plane \mathcal{A}^* is sometimes called “perpendicular regression” plane.

The values $\hat{\lambda}_i^{(n)}$, and $\hat{\alpha}_k^{(n)}$ form the basis of several rules of thumb for choosing an appropriate number k , i.e., for selecting how many principal components to retain to capture most variability in the data. Often, k is simply chosen to be the smallest number with the property $\hat{\alpha}_k^{(n)} > c$, where c is some constant, for instance 0.8. The **Kaiser's rule** suggests to take the smallest k , such that $\hat{\lambda}_k^{(n)}$ is larger than the average of all values $\hat{\lambda}_1^{(n)}, \dots, \hat{\lambda}_p^{(n)}$ ¹⁰. Another popular method is to draw the explained-variance plot which is a piece-wise linear function connecting the points $(0, 0), (1, \hat{\alpha}_1^{(n)}), \dots, (p, \hat{\alpha}_p^{(n)})$. If this line forms an “**elbow**” at a point k , it suggests that k could be an appropriate number of principal components to summarize the original data.

1.4 History and applications of principal components

The principal component analysis as a technique has been proposed in 1901 by Karl Pearson, an English biostatistician. The term itself was proposed in 1930's by Harold Hotelling, an American statistician.

The simplest PCA, as described in this chapter, is intended to be applied to a random sample of p -dimensional vectors without specific structure and without division of variables into the sets of dependent and independent variables. In this respect, PCA is called an "unsupervised method" in the machine learning and artificial intelligence community. Although some theoretical results about principal components assume normality, it is routinely applied also to non-normal data.¹¹

The results of PCA are sometimes used as an end in itself, for instance as a means for visualization of multidimensional data using the first two¹² components of the vectors of scores or for getting an insight into the data based on a possible interpretation of coefficient vectors $\hat{u}_1^{(n)}, \dots, \hat{u}_k^{(n)}$. Note, however, that there are also many other methods for the visualisation of multidimensional data, such as various projection pursuit methods, or the so-called multidimensional scaling, which we will describe in the next section.

¹⁰Note that the average of eigenvalues $\hat{\lambda}_i^{(n)}$, $i = 1, \dots, p$, is equal to the average of variances of variables X_1, \dots, X_p , which is in turn equal to $\text{tr}(S_n)/p$.

¹¹Indeed, when the sample principal components are used for visual display of data, i.e., for exploration of the structure of data, assumption of normality makes little sense. Moreover, real-world, highly multidimensional data rarely follow a multivariate normal distribution.

¹²It can be expected that in near future 3-dimensional data visualisation based on principal components will also become common.

Sometimes the results of PCA (usually the vectors of first k coordinates of scores) form an input to a subsequent statistical procedure that benefits from a small-dimensional representation of the originally large-dimensional dataset; an example is reducing the number of explanatory variables for a regression analysis. That is PCA is also used as pre-processing tool of dimensionality reduction.

PCA is used across all fields that handle multivariate data. A particularly famous application uses principal components for face recognition, see, e.g.,

<http://en.wikipedia.org/wiki/Eigenface>

but PCA is routinely applied in basically all areas that analyse multivariate numerical data, such as life sciences, genetics, chemistry, psychology, finance and econometrics, environmental science and meteorology and so on.

The most problematic aspect of principal components is its dependence on the scale of individual variables, i.e., principal components are not scale invariant. For instance, if we change the units with which we measure some distance variable (say, from meters to millimetres), the principal components can significantly change. This is particularly problematic if the variables have very different magnitudes or if we simultaneously use variables of completely different nature (such as distances, times and weights) where it is impossible to express all variables in the same units. This is the reason why PCA is sometimes based on the correlation matrix, instead of the variance-covariance matrix, that is, all variables are scaled to unit standard deviation before the application of principal components.

2 Multidimensional Scaling

Suppose that we study a set of n objects and the relations of the objects are described by an $n \times n$ matrix D of their mutual **dissimilarities**. Multidimensional scaling is a class of methods that assign vectors $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n \in \mathbb{R}^k$ to the objects, such that the mutual distances of pairs $\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j$ are “close” to dissimilarities D_{ij} of the objects i and j . The usual aim is to discover a hidden structure in the data by means of visualizing the “map” of dissimilarities in a two or three dimensional space, i.e., $k = 2$ or $k = 3$.

As we saw in the previous section, a reasonable k -dimensional representation of the data can be obtained from sample principal components, but for the direct application of principal components analysis, we need to know the

p -dimensional vectors of features of all objects. The so-called metric multidimensional scaling¹³ is closely related to PCA; here, however, only the matrix D of dissimilarities is required. Note that methods of multidimensional scaling usually do not make assumptions about the probability distribution that generated the data.

2.1 Theory for Classical Multidimensional Scaling

First, we will describe a classical solution of the following problem: How do we construct an n -tuple $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$ of points in \mathbb{R}^k , $k \leq n$, if we only know Euclidean distances between these points? Clearly, the solution is not unique, because for any solution, an orthogonal rotation or a shift also provides a feasible solution. Therefore, we can search for a solution $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$ with the center of mass in $\mathbf{0}_k$, that is, $\sum_{i=1}^n \tilde{\mathbf{x}}_i = \mathbf{0}_k$.

It is simple to prove the following lemma.

Lemma 2.1. *Let $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ be an $n \times k$ matrix, let $\sum_{i=1}^n \mathbf{x}_i = \mathbf{0}_k$, and let $B = XX^T$. Then (i) the sum of all elements of any row (and any column) of B is zero. (ii) $\|\mathbf{x}_i - \mathbf{x}_j\|^2 = B_{ii} + B_{jj} - 2B_{ij}$.*

Matrix $B = XX^T$ from the previous lemma is called the **Gram matrix** of vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ (its ij -th element is the scalar product of \mathbf{x}_i and \mathbf{x}_j). The following theorem¹⁴ shows that the Gram matrix of vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ can be computed from the mutual distances of $\mathbf{x}_1, \dots, \mathbf{x}_n$, element by element.

Theorem 2.1. *Let $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ be a matrix of the type $n \times k$, let $\sum_{i=1}^n \mathbf{x}_i = \mathbf{0}_k$ and let $B = XX^T$. Denote $D_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|$ for $i, j \in \{1, \dots, n\}$. Then*

$$B_{ij} = -\frac{1}{2} \left(D_{ij}^2 - \frac{1}{n} \sum_{r=1}^n D_{ir}^2 - \frac{1}{n} \sum_{l=1}^n D_{lj}^2 + \frac{1}{n^2} \sum_{l=1}^n \sum_{r=1}^n D_{lr}^2 \right), \quad (4)$$

for all $i, j \in \{1, \dots, n\}$.

¹³There is also an interesting variant called non-metric multidimensional scaling.

¹⁴sometimes called the Young-Householder theorem

Proof. Fix $i, j \in \{1, \dots, n\}$. Using Lemma 2.1, we have

$$\sum_{r=1}^n D_{ir}^2 = \sum_{r=1}^n (B_{ii} + B_{rr} - 2B_{ir}) = nB_{ii} + \text{tr}(B), \quad (5)$$

$$\sum_{l=1}^n D_{lj}^2 = \sum_{l=1}^n (B_{ll} + B_{jj} - 2B_{lj}) = nB_{jj} + \text{tr}(B), \quad (6)$$

$$\sum_{r=1}^n \sum_{l=1}^n D_{lr}^2 = \sum_{r=1}^n (nB_{rr} + \text{tr}(B)) = 2n\text{tr}(B). \quad (7)$$

From (7) we obtain $\text{tr}(B) = (2n)^{-1} \sum_{r=1}^n \sum_{l=1}^n D_{lr}^2$, which can be substituted to (5) and (6), yielding

$$B_{ii} = \frac{1}{n} \sum_{r=1}^n D_{ir}^2 - \frac{1}{2n^2} \sum_{r=1}^n \sum_{l=1}^n D_{lr}^2, \quad (8)$$

$$B_{jj} = \frac{1}{n} \sum_{l=1}^n D_{lj}^2 - \frac{1}{2n^2} \sum_{r=1}^n \sum_{l=1}^n D_{lr}^2. \quad (9)$$

From Lemma 2.1 we see that $B_{ij} = -\frac{1}{2} (D_{ij}^2 - B_{ii} - B_{jj})$, which together with (8) and (9) provides the equality from the statement of the theorem. \square

Exercise 2.1. Show that the matrix B from the previous theorem can be obtained using the following matrix computation. Let $P = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$ ¹⁵ and let A be the matrix with elements $A_{ij} = -\frac{1}{2} D_{ij}^2$, $i, j \in \{1, \dots, n\}$. Then $B = PAP$.

Therefore, distances of vectors directly provide the Gram matrix of the vectors. We will show that from the Gram matrix it is a simple step to obtain a solution $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$ of the original problem. In other words, if $B = XX^T$ for some $n \times k$ matrix X , then we can easily find some $n \times p$ matrix \tilde{X} , such that $B = \tilde{X}\tilde{X}^T$.

Clearly, the Gram matrix $B = XX^T$ is non-negative definite, that is, $B = U\Lambda U^T$, where $U = (\mathbf{u}_1, \dots, \mathbf{u}_n)$ is an orthogonal matrix of eigenvectors of B and Λ is a diagonal matrix with eigenvalues $\lambda_1 \geq \dots \geq \lambda_n \geq 0$ on the diagonal. Since X is of the type $n \times k$, where $k \leq n$, the rank of B is at most k , i.e., $\lambda_{k+1} = \dots = \lambda_n = 0$. Now we can easily verify that the matrix $\tilde{X} = (\sqrt{\lambda_1} \mathbf{u}_1, \dots, \sqrt{\lambda_k} \mathbf{u}_k)$ satisfies $B = \tilde{X}\tilde{X}^T$. That is, the k -dimensional columns $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$ of \tilde{X}^T are the required solutions.

¹⁵Note that P is a projector.

It is interesting to note that as a method of “dimensionality reduction”, multidimensional scaling is essentially equivalent to the principal components analysis. In other words, if we do have p -dimensional vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ of features of n objects (or measurements of p variables on n objects), we can decide to compute the matrix D of mutual Euclidean distances of $\mathbf{x}_1, \dots, \mathbf{x}_n$, and then use the classical multidimensional scaling with some $k \leq p$, as described above. Then, the resulting vectors $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n \in \mathbb{R}^k$ are just orthogonally rotated vectors of the first k -coordinates of the principal component scores.

2.2 Application of Classical Multidimensional Scaling

For real objects, the matrix D of dissimilarities is usually not based on Euclidean distances of some unknown vectors. Nevertheless, if we tentatively assume that D is a matrix of Euclidean distances of some vectors in \mathbb{R}^k (or in a k -dimensional subspace of \mathbb{R}^p) and imitate the theoretical construction of the vectors $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$ from the previous subsection, we often obtain a reasonable fit with the dissimilarities given by D .

If we use the matrix D of dissimilarities, it is always possible to compute the symmetric matrix B with elements given in Theorem 2.1. If D is not a perfect matrix of Euclidean distances, then B is not necessarily non-negative definite. However, if k largest eigenvalues $\lambda_1, \dots, \lambda_k$ of B are positive and large compared to the absolute values of the remaining eigenvalues¹⁶, then the mutual Euclidean distances of the k -dimensional rows of $(\sqrt{\lambda_1}\mathbf{u}_1, \dots, \sqrt{\lambda_k}\mathbf{u}_k)$ give a good fit with dissimilarities D_{ij} .

Classical multidimensional scaling, as a means to visualize data given by their matrix of dissimilarities, have been used for instance in psychology to create a “perceptual map” of stimuli, in marketing to create a “product (diss)similarity map”, in social networks to create “friendship/collaboration maps” and so on.

Note that there are several more advanced alternatives to the classical multidimensional scaling: the so-called general metric multidimensional scaling and non-metric multidimensional scaling. These methods are beyond the scope of this introductory lecture.

¹⁶A common rule of thumb is that $\sum_{i=1}^k \lambda_i / \sum_{i=1}^n |\lambda_i| > 0.8$.

3 Canonical Correlations

3.1 Mathematical background for Canonical Correlations

The theory of canonical correlation can be easily explained using the notion of the square-root matrix. Let Σ be a non-negative definite $p \times p$ matrix and let $\Sigma = \sum_{i=1}^p \lambda_i u_i u_i^T$ be the decomposition of Σ from Theorem 1.1. Let $\Sigma^{1/2} := \sum_{i=1}^p \sqrt{\lambda_i} u_i u_i^T$. Clearly, $\Sigma^{1/2}$ is a non-negative definite matrix satisfying $\Sigma^{1/2} \Sigma^{1/2} = \Sigma$, i.e., it is natural to call it the **square-root matrix** of Σ ¹⁷. If Σ is positive definite, its square-root matrix is also positive definite and its inverse $\Sigma^{-1/2} := (\Sigma^{1/2})^{-1}$ satisfies $\Sigma^{-1/2} = \sum_{i=1}^p \lambda_i^{-1/2} u_i u_i^T$.

Let $u \in \mathbb{R}^q$ be a vector of norm 1 and let M be a non-negative definite matrix of the type $q \times q$. Then $u^T M u \leq \lambda_1(M)$, where $\lambda_1(M)$ is the largest eigenvalue of M , which is sometimes called the **Railegh-Ritz theorem**. Similarly, assume that M has $m \leq p$ distinct positive eigenvalues $\lambda_1(M) > \dots > \lambda_p(M)$ (and $p - m$ zero eigenvalues). Let $2 \leq k \leq m$. If u is orthogonal on $k - 1$ eigenvectors of M corresponding to the $k - 1$ largest eigenvalues of M , then $u^T M u \leq \lambda_k(M)$.

Recall also that, if F is any matrix, then the linear space generated by the columns of F is the same as the linear space generated by the columns of the matrix FF^T . If A, B are $s \times r$ matrices then $\text{tr}(A^T B) = \text{tr}(BA^T) = \text{tr}(AB^T) = \text{tr}(B^T A)$. Note that the largest eigenvalue of a non-negative definite matrix of rank 1 is equal to the trace of the matrix.

3.2 Theoretical Canonical Correlations

Consider the random vector $\mathbf{X} = (\mathbf{X}_{(1)}^T, \mathbf{X}_{(2)}^T)^T$ where $\mathbf{X}_{(1)}$ is the subvector of dimension p and $\mathbf{X}_{(2)}$ is the subvector of dimension q . Correspondingly, let the covariance matrix of \mathbf{X} be divided into sub-blocks as follows:

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}.$$

¹⁷It is also possible to show that $\Sigma^{1/2}$ is unique even in the case when the orthonormal system u_1, \dots, u_p of eigenvectors of Σ is not uniquely defined.

We will assume that Σ_{11} and Σ_{22} are positive definite, i.e., there exist matrices $\Sigma_{11}^{-1/2}$ and $\Sigma_{22}^{-1/2}$. Let

$$\begin{aligned} B &:= \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2}, \\ N_1 &:= BB^T = \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2}, \\ N_2 &:= B^T B = \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1/2}. \end{aligned}$$

Clearly, N_1, N_2 are both non-negative definite. Moreover, N_1, N_2 have the same rank, because $\text{rank}(N_1) = \text{rank}(BB^T) = \text{rank}(B) = \text{rank}(B^T) = \text{rank}(B^T B) = \text{rank}(N_2)$. We will denote the common rank of matrices N_1, N_2 by the symbol m .

Lemma 3.1. *Let $\alpha_1, \dots, \alpha_m$ be an orthonormal system of eigenvectors of N_1 , with corresponding eigenvalues $\lambda_1(N_1) > \dots > \lambda_m(N_1) > 0$. Similarly, let β_1, \dots, β_m be the orthonormal system of eigenvectors of N_2 , with corresponding eigenvalues $\lambda_1(N_2) > \dots > \lambda_m(N_2) > 0$. Then, for every $i = 1, \dots, m$, we have $\lambda_i(N_1) = \lambda_i(N_2) =: \lambda_i$ and*

$$\beta_i = \frac{B^T \alpha_i}{\sqrt{\lambda_i}}, \quad \text{or} \quad \beta_i = -\frac{B^T \alpha_i}{\sqrt{\lambda_i}}.$$

Proof. For the proof that all non-zero eigenvalues of N_1 and N_2 are equal, it is enough to show that each eigenvalue of N_1 is an eigenvalue of N_2 ¹⁸. For any $i \in \{1, \dots, m\}$, a normalized eigenvector α_i of $N_1 = BB^T$ and corresponding eigenvalue $\lambda_i(N_1)$ we can write:

$$\begin{aligned} BB^T \alpha_i &= \lambda_i(N_1) \alpha_i, \\ B^T BB^T \alpha_i &= \lambda_i(N_1) B^T \alpha_i, \\ N_2 B^T \alpha_i &= \lambda_i(N_1) B^T \alpha_i \end{aligned} \tag{10}$$

In addition, observe that

$$\|B^T \alpha_i\|^2 = \alpha_i^T BB^T \alpha_i = \alpha_i^T N_1 \alpha_i = \alpha_i^T \lambda_i(N_1) \alpha_i = \lambda_i(N_1) > 0, \tag{11}$$

which implies $B^T \alpha_i \neq 0$. That is (10) entails that $\lambda_i(N_1)$ is an eigenvalue of N_2 and $B^T \alpha_i$ is a corresponding eigenvector of N_2 (not necessarily normalised). Moreover, since all eigenvalues of N_2 are assumed to be distinct, the corresponding orthonormal system of eigenvectors of N_2 is uniquely determined, up to the possible reversal of the direction. Therefore

$$\beta_i = \pm \frac{B^T \alpha_i}{\|B^T \alpha_i\|} = \pm \frac{B^T \alpha_i}{\sqrt{\lambda_i}},$$

¹⁸Clearly, then, the symmetry of the problem implies that each eigenvalue of N_2 will be an eigenvalue of N_1 .

where the second equality follows from (11). □

Definition 3.1 (Canonical variables and canonical correlations). For $i = 1, \dots, m$ let

$$\begin{aligned} a_i &= \Sigma_{11}^{-1/2} \alpha_i, \quad b_i = \Sigma_{22}^{-1/2} \beta_i, \\ U_i &= a_i^T \mathbf{X}_{(1)}, \quad V_i = b_i^T \mathbf{X}_{(2)}, \\ \rho_i &= \sqrt{\lambda_i}, \end{aligned}$$

where α_i, β_i and λ_i are defined in Lemma 3.1, such that $\beta_i = \frac{\mathbf{B}^T \alpha_i}{\sqrt{\lambda_i}}$. Then random variables $U_1, \dots, U_m, V_1, \dots, V_m$ are called **canonical variables** and numbers $\rho_1, \dots, \rho_m, 0, \dots, 0$ ¹⁹ are called **canonical correlations**.

It is possible to show that the vectors a_i of coefficients that define the first group of canonical variables are eigenvectors of $\Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$, and canonical correlations ρ_i are square roots of corresponding eigenvalues. Analogously, vectors b_i of coefficients that define the second group of canonical variables are eigenvectors of $\Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$ and, as above, canonical correlations ρ_i are square roots of corresponding eigenvalues. That is, we can obtain canonical variables and correlations without actually computing the square root matrices; however, the square root matrices are useful for the theoretical results, because they allow us to work exclusively with non-negative definite matrices.

Theorem 3.1 (Mutual correlations of canonical variables). For all $i, j \in \{1, \dots, m\}$ we have

$$\begin{aligned} \text{cov}(U_i, U_j) &= \rho(U_i, U_j) = \delta_{ij}, \\ \text{cov}(V_i, V_j) &= \rho(V_i, V_j) = \delta_{ij}, \\ \text{cov}(U_i, V_j) &= \rho(U_i, V_j) = \delta_{ij} \rho_i, \end{aligned}$$

where δ_{ij} is the Kronecker delta²⁰.

Proof.

$$\begin{aligned} \text{cov}(U_i, U_j) &= \text{cov}(a_i^T \mathbf{X}_{(1)}, a_j^T \mathbf{X}_{(1)}) = a_i^T \Sigma_{11} a_j \\ &= \alpha_i^T \Sigma_{11}^{-1/2} \Sigma_{11} \Sigma_{11}^{-1/2} \alpha_j = \alpha_i^T \alpha_j = \delta_{ij} \end{aligned}$$

¹⁹The number of zeros here is $\min(p, q) - m$, that is all canonical correlations of order higher than $m = \text{rank}(N_1) = \text{rank}(N_2)$ are defined to be zero.

²⁰ $\delta_{ii} = 1$ and $\delta_{ij} = 0$ if $i \neq j$.

and, similarly, we obtain $\text{cov}(V_i, V_j) = \delta_{ij}$. In particular, this implies $D(U_i) = D(V_i) = 1$, that is, the covariances are equal to correlations. We can conclude the proof by observing that:

$$\begin{aligned} \text{cov}(U_i, V_j) &= \text{cov}(a_i^T \mathbf{X}_{(1)}, b_j^T \mathbf{X}_{(2)}) = a_i^T \Sigma_{12} b_j = \alpha_i^T \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2} \beta_j \\ &= \alpha_i^T \mathbf{B} \beta_j =^* \sqrt{\lambda_i} \beta_i^T \beta_j = \delta_{ij} \rho_i, \end{aligned}$$

where for the equality denoted by the asterisk we used $\beta_i = \frac{\mathbf{B}^T \alpha_i}{\sqrt{\lambda_i}}$. \square

In a manner similar to principal components, canonical variables and canonical correlations have a probabilistic justification:

Theorem 3.2 (Maximum correlation justification of canonical variables and correlations). *Canonical variables $U_1 = a_1^T \mathbf{X}_{(1)}$ and $V_1 = b_1^T \mathbf{X}_{(2)}$ have the maximum possible correlation among all pairs $a^T \mathbf{X}_{(1)}$, $b^T \mathbf{X}_{(2)}$ of non-zero linear combinations of the components of vectors $\mathbf{X}_{(1)}$ and $\mathbf{X}_{(2)}$. Formally:*

$$\rho(U_1, V_1) \geq \rho(a^T \mathbf{X}_{(1)}, b^T \mathbf{X}_{(2)}) \text{ for all } 0 \neq a \in \mathbb{R}^p, 0 \neq b \in \mathbb{R}^q.$$

For $k \geq 2$, the canonical correlations $U_k = a_k^T \mathbf{X}_{(1)}$ and $V_k = b_k^T \mathbf{X}_{(2)}$ have the largest correlation coefficient among all pairs $a^T \mathbf{X}_{(1)}$, $b^T \mathbf{X}_{(2)}$ of non-zero linear combinations of the components of vectors $\mathbf{X}_{(1)}$ and $\mathbf{X}_{(2)}$ that are uncorrelated with U_1, \dots, U_{k-1} and V_1, \dots, V_{k-1} , respectively. Formally:

$$\begin{aligned} \rho(U_k, V_k) &\geq \rho(a^T \mathbf{X}_{(1)}, b^T \mathbf{X}_{(2)}) \text{ for all } 0 \neq a \in \mathbb{R}^p, 0 \neq b \in \mathbb{R}^q \\ \text{s.t. } a^T \Sigma_{11} a_i &= b^T \Sigma_{22} b_i = 0 \text{ for all } i = 1, \dots, k-1. \end{aligned}$$

Proof. As $\rho(U_1, V_1) = \sqrt{\lambda_1}$ and

$$\rho(a^T \mathbf{X}_{(1)}, b^T \mathbf{X}_{(2)}) = \frac{a^T \Sigma_{12} b}{\sqrt{a^T \Sigma_{11} a} \sqrt{b^T \Sigma_{22} b}},$$

which is invariant with respect to multiples of vectors a, b , for the proof of the first part of the theorem it is enough to show that $\lambda_1 \geq (a^T \Sigma_{12} b)^2$ under the constraint $a^T \Sigma_{11} a = b^T \Sigma_{22} b = 1$. Denote $u := \Sigma_{11}^{1/2} a$ and $v := \Sigma_{22}^{1/2} b$ and note that $a^T \Sigma_{11} a = b^T \Sigma_{22} b = 1$ implies $\|u\| = \|v\| = 1$. We obtain:

$$\begin{aligned} (a^T \Sigma_{12} b)^2 &= (u^T \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2} v)^2 \\ &= u^T \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2} v v^T \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1/2} u \\ &\leq^* \lambda_1 (\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2} v v^T \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1/2}) \\ &= \text{tr}(\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2} v v^T \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1/2}) \\ &= v^T \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1/2} \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2} v \\ &= v^T N_2 v \leq^* \lambda_1(N_2) = \lambda_1. \end{aligned}$$

The inequalities denoted by the asterisk follow from the Raileigh-Ritz inequality.

The second part of the theorem follows similarly as the first one, observing that the conditions $b^T \Sigma_{22} b_i = 0$ for all $i = 1, \dots, k - 1$ imply that v is orthogonal to the eigenvectors of N_2 that correspond to $k - 1$ largest eigenvalues, hence $v^T N_2 v \leq \lambda_k(N_2) = \lambda_k$. \square

Note that unlike principal components, canonical correlations ρ_1, \dots, ρ_m do not depend on the units of individual variables, i.e., they are scale invariant.

3.3 Sample Canonical Correlations

In practice, we do not know the theoretical variance-covariance matrices of the random vector that we observe and we need to based the estimate of the canonical correlations on a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ from the underlying $(p+q)$ -dimensional distribution. The following definition is motivated by the remark after Definition 3.1.

Definition 3.2 (Sample canonical correlations). *Let $m = \min(p, q)$. Sample canonical correlations $\hat{\rho}_1, \dots, \hat{\rho}_m$ are defined to be the square roots of the m largest eigenvalues of $\mathbf{S}_{22}^{-1} \mathbf{S}_{21} \mathbf{S}_{11}^{-1} \mathbf{S}_{12}$ (or, equivalently, of $\mathbf{S}_{11}^{-1} \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{21}$), where $\mathbf{S}_{11}, \mathbf{S}_{12}, \mathbf{S}_{21}$, and \mathbf{S}_{22} are the $p \times p$, $p \times q$, $q \times p$, and $q \times q$ sub-blocks of the sample variance-covariance matrix*

$$\mathbf{S} = \begin{pmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{21} & \mathbf{S}_{22} \end{pmatrix}.$$

In the case of normal errors and a large sample size, we can test the hypothesis that all theoretical canonical correlations are zero, i.e., that the sub-vector of the first p variables is independent with the sub-vector of the last q variables:

Theorem 3.3. *If $\mathbf{X}_1, \dots, \mathbf{X}_n$ follow $N_{p+q}(\mu, \Sigma)$ and the upper-right $p \times q$ submatrix Σ_{12} of Σ is zero²¹, then the statistics*

$$W = \det(\mathbf{I} - \mathbf{S}_{22}^{-1} \mathbf{S}_{21} \mathbf{S}_{11}^{-1} \mathbf{S}_{12})$$

has asymptotically the Wilks distribution $\Lambda(p, n - 1 - q, q)$ and the statistics

$$Z = -(n - (p + q + 3)/2) \ln(W)$$

has asymptotically the distribution χ_{pq}^2 .

The tests using the Z statistics from the previous theorem is called the Bartlett ξ^2 test.

²¹This is of course if and only if the lower-left $q \times p$ submatrix Σ_{21} of Σ is zero.

3.4 Applications of Canonical Correlations

Canonical correlations are used in the situations where variables can be logically divided into two distinct groups. For instance, in a psychological survey we can have a set of variables measuring one “personality dimension”, and a set of variables measuring a second “personality dimension”. The first canonical correlation ρ_1 then measures the overall degree of correlation of the first group of variables with the second group of variables.

Note that for $q = 1$ the first canonical correlation corresponds to the so-called coefficient of multiple correlation, which, in the context of linear regression, is closely related to the square of the coefficient of determination. (It is a useful theoretical exercise to simplify the theory of canonical correlations for $q = 1$.)

4 Factor Analysis

In factor analysis, we assume that the random vector $\mathbf{X} = (X_1, \dots, X_p)^T$ of observed (or “manifest”) variables can be described in terms of an m -dimensional random vector $\mathbf{F} = (F_1, \dots, F_m)^T$ of hidden (or “latent”) variables, where m is substantially smaller than p . More precisely, we assume that the following statistical model holds:

$$X_i = \sum_{j=1}^m a_{ij}F_j + U_i, \quad i = 1, \dots, p, \quad (12)$$

or, in a matrix form,

$$\mathbf{X} = \mathbf{A}\mathbf{F} + \mathbf{U}. \quad (13)$$

The elements a_{ij} of \mathbf{A} are called **factor loadings**²², the random vector \mathbf{F} is called the vector of **common factors**, and the random vector \mathbf{U} is called the vector of **specific factors**, specific variates, or uniqueness. The factor loadings are assumed to be unknown parameters of the model.

Hence, we consider a model similar to the multivariate regression, but in factor analysis the “regressors” or “explanatory variables” F_1, \dots, F_m are assumed to be random, and cannot be directly observed²³.

²²That is, the $p \times m$ matrix \mathbf{A} is called the matrix of factor loadings.

²³Another difference is that in factor analysis, we allow the variances of the components of the “error vector” \mathbf{U} to differ, while in regression analysis the errors are usually assumed to be homoscedastic.

Usual theoretical assumptions are that the common factors \mathbf{F} and uniqueness \mathbf{U} are uncorrelated, i.e., $Cov(\mathbf{F}, \mathbf{U}) = 0_{m \times p}$, the common factors themselves are standardized and uncorrelated, which means that the variance-covariance matrix of \mathbf{F} is I_m ²⁴, and the variance-covariance matrix of \mathbf{U} is assumed to be a diagonal matrix $D = \text{diag}(d_1, \dots, d_p)$. Note that the variances d_1, \dots, d_p are additional unknown parameters of the model.

The mean values of \mathbf{F} and \mathbf{U} are assumed to be 0, which directly implies that the mean value of \mathbf{X} is 0. In this text, we also adopt a simplifying assumption that the manifest variables X_1, \dots, X_p are normalized, i.e., $D(X_i) = 1$, which is common in applications²⁵.

The theoretical assumptions, together with the standard transformation rules stated in Theorem 1.2 imply that the correlation matrix Ψ of the vector \mathbf{X} ²⁶ is

$$\Psi = AA^T + D. \quad (14)$$

Note that the elements of the matrix A are simply correlations between the manifest variables and hidden factors:

Theorem 4.1 (Interpretation of factor loadings). *Under the model of factor analysis described above, we have $a_{ij} = \rho(X_i, F_j)$ for all $i = 1, \dots, p$ and $j = 1, \dots, m$.*

Proof. Consider the model (13). We have $Cov(\mathbf{X}, \mathbf{F}) = Cov(A\mathbf{F} + \mathbf{U}, \mathbf{F}) = ACov(\mathbf{F}, \mathbf{F}) + Cov(\mathbf{U}, \mathbf{F}) = A$, taking into account assumptions $Cov(\mathbf{F}, \mathbf{F}) = I_m$ and $Cov(\mathbf{U}, \mathbf{F}) = 0_{p \times m}$. Since the variables of \mathbf{X} as well as of \mathbf{F} are normalized, we obtain that $\rho(X_i, F_j) = cov(X_i, F_j) = a_{ij}$. \square

A crucial observation is that the model (13) is over-parametrized: Based on the observations of the manifest variables \mathbf{X} , we will never be able to distinguish between models $\mathbf{X} = A\mathbf{F} + \mathbf{U}$ and $\mathbf{X} = (AV)(V^T\mathbf{F}) + \mathbf{U}$, where V is any $m \times m$ orthogonal matrix. In other words, the factor loadings A with factors \mathbf{F} provide as good a fit to the observations as factor loadings AV with “rotated” factors $V^T\mathbf{F}$ ²⁷.

²⁴This is then called the orthogonal model, as opposed to a non-orthogonal model, where no restrictions on the variance-covariance matrix of \mathbf{F} are imposed.

²⁵Especially in some behavioral sciences, where X_i represent standardized responses of subjects.

²⁶Note that since the manifest variables are standardized, the correlation matrix of \mathbf{X} coincides with the variance-covariance matrix Σ of \mathbf{X} .

²⁷Note that if \mathbf{F} satisfies the assumptions of factor analysis, then so does $V^T\mathbf{F}$.

The previous point is the crux of the main criticism of factor analysis from the point of view of traditional statistical philosophy. Note that statisticians traditionally assume a single, fixed model, which is unknown, yet possible to identify with increasing precision as the sample size increases. Clearly, the model of factor analysis does not fit into this view. Rather, in factor analysis, we can adopt the attitude that *all* models of the form $\mathbf{X} = (AV)(V^T \mathbf{F}) + \mathbf{U}$, where V is an orthogonal matrix, are equally good representations of reality, and it is perfectly justifiable to select one of these models that leads to our understanding of the origin of data, based on a clear interpretation of what do the factors $V^T \mathbf{F}$ represent. In fact, there is a large body of literature on methods called “rotation of factors” that help us find a suitable matrix V ; see the section on rotation of factors.

4.1 Estimation of factor loadings

The first step in factor analysis is to use the observations of the manifest variables to find some appropriate estimates of the parameters a_{ij} and d_i . Here, two methods are most popular: **the method of maximal likelihood** and **the method of principal factors** which we very briefly describe.

Consider the formula $\Psi = AA^T + D$. If we knew $D = \text{diag}(d_1, \dots, d_p)$, then we would be able to evaluate the so-called **reduced correlation matrix** $\Psi - D = AA^T$, which would enable us the computation of a matrix A of factor loadings. Note that $d_i = 1 - h_i^2$, where $h_i^2 = \sum_{j=1}^m a_{ij}^2$ are called **communalities**, i.e., the estimation of the variances of specific factors is essentially equivalent to the estimation of communalities.

The method of principal factors uses the sample correlation matrix R as an estimator of Ψ , and $\hat{d}_1 = 1/(R^{-1})_{11}, \dots, \hat{d}_p = 1/(R^{-1})_{pp}$ as (initial) estimators of d_1, \dots, d_p . Thus, we can obtain an estimate

$$R^* = R - \text{diag}(\hat{d}_1, \dots, \hat{d}_p)$$

of the reduced correlation matrix, which can be viewed as the sample correlation matrix with diagonal elements replaced by the estimates $\hat{h}_1^2 = 1 - 1/(R^{-1})_{11}, \dots, \hat{h}_p^2 = 1 - 1/(R^{-1})_{pp}$ of the communalities. Then we calculate the eigenvalues $\lambda_1 \geq \dots \geq \lambda_p$ of R^* and the corresponding orthonormal eigenvectors u_1, \dots, u_p . If the first m eigenvalues of R^* are non-negative, and the remaining eigenvalues are “small”, we can estimate the matrix of factor loadings as

$$\hat{A} = (\sqrt{\lambda_1}u_1, \dots, \sqrt{\lambda_m}u_m).$$

We remark that once we have \hat{A} , we can repeat the process using new estimates $\sum_{j=1}^m (\hat{A})_{ij}^2$, $i = 1, \dots, p$, of communalities and iterate it until possible convergence.

In the method of principal factors the number m of factors is considered satisfactory based on analogous criteria as in the method of principal components, for instance, if $\sum_{j=1}^m \lambda_j / \sum_{i=1}^p |\lambda_i| > 0.8$.

4.2 Rotation of Factors

As mentioned above, if we have any estimate \hat{A} of A , and V is any $m \times m$ orthogonal matrix, then $\tilde{A} = \hat{A}V$ is as good an estimate of A as \hat{A} . For a given estimate \hat{A} , the methods of factors rotation provide, in a sense, optimal rotation V^* in order to achieve a suitable form of the new matrix $\hat{A}V^*$ of factor loadings. A general aim is to achieve the so called **simple structure** of the matrix of factor loadings. Omitting details, a matrix of factor loadings has a simple structure, if it contains many entries close to 0, 1 and -1 , because, viewed as correlations, such factor loadings usually lead to a simple interpretation of common factors.

Definition 4.1 (Varimax and Quartimax rotation methods). *Varimax rotation of the $p \times m$ matrix \hat{A} of factor loadings is the orthogonal matrix V of the type $m \times m$ that maximizes the value of*

$$f_v(V) = \frac{1}{m} \sum_{t=1}^m \left[\frac{1}{p} \sum_{j=1}^p \left((\hat{A}V)_{jt}^2 - \frac{1}{p} \sum_{k=1}^p (\hat{A}V)_{kt}^2 \right)^2 \right]. \quad (15)$$

Quartimax rotation of the $p \times m$ matrix \hat{A} of factor loadings is the orthogonal matrix V of the type $m \times m$ that maximizes the value of

$$f_q(V) = \frac{1}{mp} \sum_{t=1}^m \sum_{j=1}^p \left((\hat{A}V)_{jt}^2 - \frac{1}{mp} \sum_{s=1}^m \sum_{k=1}^p (\hat{A}V)_{ks}^2 \right)^2. \quad (16)$$

Let B_V be the $m \times p$ matrix with elements corresponding to the squares of the elements of $\hat{A}V$, that is, $B_V = \hat{A}V \odot \hat{A}V$, where \odot is the Hadamard (entry-wise) product of matrices. Observe that it is possible to interpret (15) as the average “sample variance” of the entries of the columns of B_V , and (16) can be interpreted as the “sample variance” of all elements of B_V . Clearly, the elements of $\hat{A}V$ are correlations, i.e., the elements of B_V are bounded by 0 from below and by 1 from above. Therefore, the maximization of the

“variance” of the elements of B_V forces the elements to be close to 0 or 1, i.e., forces the elements of $\hat{A}V$ to be close to the numbers 0, 1, and -1 .

Using the formula $\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \frac{1}{N} \sum_{i=1}^N x_i^2 - \bar{x}^2$ for any x_1, \dots, x_N and the observation

$$\sum_{s=1}^m \sum_{k=1}^p (\hat{A}V)_{ks}^2 = \|\hat{A}V\|^2 = \text{tr}(\hat{A}V(\hat{A}V)^T) = \text{tr}(\hat{A}\hat{A}^T) = \|\hat{A}\|^2,$$

the quartimax utility function can be simplified as follows:

$$f_q(V) = \frac{1}{mp} \sum_{t=1}^m \sum_{j=1}^p (\hat{A}V)_{jt}^4 - \frac{1}{(mp)^2} \|\hat{A}\|^4.$$

Therefore, the quartimax method maximizes the sum of the fourth powers of factor loadings, hence the name.

The problem of finding optimal V in the sense of the varimax or the quartimax rotation method is a difficult problem of mathematical programming and description of the appropriate numerical methods goes beyond this text. However, all advanced statistical packages implement some form of factor rotation optimization.

4.3 Estimation of Factor Scores

Lemma 4.1. *Let the random vector $(\mathbf{X}^T, \mathbf{F}^T)^T$ follows the $p+m$ -dimensional normal distribution with the zero mean value. Then the conditional distribution of \mathbf{F} given $\mathbf{X} = \mathbf{x}$ is*

$$N_m \left(\text{Cov}(\mathbf{F}, \mathbf{X}) \Sigma_{\mathbf{X}}^{-1} \mathbf{x}, \Sigma_{\mathbf{F}} - \text{Cov}(\mathbf{F}, \mathbf{X}) \Sigma_{\mathbf{X}}^{-1} \text{Cov}(\mathbf{X}, \mathbf{F}) \right) = N_m \left(A^T (AA^T + D)^{-1} \mathbf{x}, I_m - A^T (AA^T + D)^{-1} A \right).$$

Consider a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ of the p -dimensional observable variables on n objects, satisfying the model of factor analysis. Suppose that our aim is to estimate the realizations of the common factors $\mathbf{F}_1, \dots, \mathbf{F}_n$ for individual objects; these realizations are called **factor scores**. Lemma 4.1 motivates the following estimators:

$$\hat{\mathbf{F}}_r = \hat{A}^T (\hat{A}\hat{A}^T + \hat{D})^{-1} \mathbf{X}_r, \quad r = 1, \dots, n, \quad (17)$$

where \hat{A} is the selected estimate of the matrix of factor loadings and \hat{D} is the estimate of the diagonal matrix of variances of specific variates. This kind of

estimation of factor scores is sometimes called the **method of regression analysis**²⁸.

If the number p is very large, it may be difficult to compute the inverse of the $p \times p$ correlation matrix $\hat{R} = \hat{A}\hat{A}^T + \hat{D}$ ²⁹. However, the matrix $\hat{A}^T\hat{R}^{-1}$ that appears in the formula (17) can be computed by only inverting a usually much smaller $m \times m$ matrix, and a $p \times p$ diagonal matrix:

Lemma 4.2. *Let $R = AA^T + D$ be a non-singular matrix, where A is a $p \times m$ matrix and D is a $p \times p$ matrix. Then $A^T R^{-1} = (I_m + A^T D^{-1} A)^{-1} A^T D^{-1}$.*

Proof.

$$\begin{aligned} (I_m + A^T D^{-1} A)^{-1} A^T D^{-1} R &= \\ (I_m + A^T D^{-1} A)^{-1} A^T D^{-1} (AA^T + D) &= \\ (I_m + A^T D^{-1} A)^{-1} (A^T D^{-1} A + I_m) A^T &= A^T. \end{aligned}$$

□

5 Partitioning Methods of Cluster Analysis

The general aim of clustering is to reveal relationships or similarities of n objects, typically characterized by either multidimensional real-valued vectors of features or by a matrix of mutual dissimilarities. In particular, the aim of the *partitioning* methods of cluster analysis³⁰ is to divide the objects into k groups called clusters, in such a way that the elements within the same cluster are as similar as possible and the elements of different clusters are as dissimilar as possible. It is the specification of the notions of “as similar as possible” and “as dissimilar as possible” that distinguishes different partitioning methods.

First, let us introduce some notation. Let n be the number of objects and let k be the required number of clusters.³¹ That is, the objects will be

²⁸The reason is that it can also be motivated by techniques similar to multivariate regression analysis.

²⁹It was especially difficult in the past.

³⁰Partitioning methods are sometimes called non-hierarchical methods to highlight the contrast with the hierarchical methods of cluster analysis.

³¹We begin with an assumption that the number of clusters is given. Of course, in many applications we do not know an appropriate number of clusters in advance. We discuss methods for choosing the correct number of clusters for such applications later in this chapter.

denoted by the numbers $1, \dots, n$ and the clusters by the numbers $1, \dots, k$. We will assume that $n \geq k$; the case $n < k$ is meaningless; the number n of objects is usually greater than k by several orders of magnitude.

Definition 5.1 (Labelling and clusters). *Let $\Gamma_{n,k}$ be the set of all vectors $\gamma \in \{1, \dots, k\}^n$, such that each of the values $1, \dots, k$ occurs at least once in γ ³². Any vector $\gamma \in \Gamma_{n,k}$ will be called a “labelling” (or a “cluster assignment”). The set $C_j(\gamma) = \{i \in \{1, \dots, n\} : \gamma_i = j\}$ will be called the j -th cluster for the labelling γ ³³. The number $n_j(\gamma)$ of elements of $C_j(\gamma)$, i.e., $n_j(\gamma) = |C_j(\gamma)|$, will be called the size of the j -th cluster for the labelling γ .*

Note that different labellings that differ only by a permutation of labels represent the same “unlabelled” partition. Most clustering methods treat two labellings γ_1 and γ_2 as equally good if they differ only by a permutation of cluster labels, i.e., if there exists a permutation π of $\{1, \dots, k\}$ such that $\gamma_1(i) = \pi(\gamma_2(i))$ for all $i \in \{1, \dots, n\}$.

If the objects are characterized by features $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$,³⁴ we can define the following characteristics for each cluster $C_j(\gamma)$, $\gamma \in \Gamma_{n,k}$.

Definition 5.2 (Centroid of a cluster). *Let $\gamma \in \Gamma_{n,k}$ and let $j \in \{1, \dots, k\}$. The centroid of the cluster $C_j(\gamma)$ is³⁵*

$$\bar{\mathbf{x}}_j(\gamma) = n_j^{-1}(\gamma) \sum_{i \in C_j(\gamma)} \mathbf{x}_i.$$

Definition 5.3 (Variance-covariance matrix of a cluster). *Let $\gamma \in \Gamma_{n,k}$ and let $j \in \{1, \dots, k\}$. The variance-covariance matrix of the cluster $C_j(\gamma)$ is³⁶*

$$S_j(\gamma) = n_j^{-1}(\gamma) \sum_{i \in C_j(\gamma)} (\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma))(\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma))^T.$$

³²As we will see, the interpretation of this requirement is that all clusters contain at least one object.

³³That is, for $j \in \{1, \dots, k\}$ and $i \in \{1, \dots, n\}$ the equality $\gamma_i = j$ means that the labelling γ assigns the i -th object into the j -th cluster.

³⁴At the moment we view $\mathbf{x}_1, \dots, \mathbf{x}_n$ just as a collection of (known) observed vectors. This is a typical framework of cluster analysis. Later, we will reframe $\mathbf{x}_1, \dots, \mathbf{x}_n$ as samples from an underlying distribution; the statistical point of view suggests some powerful clustering methods such as the Gaussian model-based clustering.

³⁵Note that the centroid of a cluster is its “center of mass”, if the objects are assumed to have mass 1 and \mathbf{x}_i is the position of the i -th object, $i = 1, \dots, n$.

³⁶The Variance-covariance matrix is a matrix representation of the “shape” of the cluster. Note, however, that in some pathological cases it can be singular. Note also that some texts use the normalization constant $(n_j(\gamma) - 1)^{-1}$ motivated by unbiased estimation in an underlying statistical model. The practical and theoretical difference is typically minor.

Partitioning methods of cluster analysis compute a labelling $\gamma^* \in \Gamma_{n,k}$,³⁷ which can, however, be viewed in different ways. A more computer-science-oriented point of view is that γ^* is the result of a well-specified algorithm, e.g., the k -means clustering discussed in the next section is simply the result of a " k -means algorithm". On the other hand, an applied-mathematical point of view is that we first choose a loss function $L(\gamma)$ over all labellings $\gamma \in \Gamma_{n,k}$, and the labelling γ^* is the solution (or an approximation of the solution) of the associated discrete optimization problem. The loss function can be based on geometric or probabilistic/statistical considerations.

5.1 Clustering using k -means and k -medoids

The best known and popular partitioning method is called " k -means". The basic idea of k -means is to first fix k and then select $\gamma_k^* \in \Gamma_{n,k}$ that minimizes

$$L_k(\gamma) := \sum_{j=1}^k \sum_{i \in C_j(\gamma)} d^2(\mathbf{x}_i, \bar{\mathbf{x}}_j(\gamma)), \quad (18)$$

where d is some distance measure, usually the Euclidean distance $d(\mathbf{x}, \tilde{\mathbf{x}}) = \|\mathbf{x} - \tilde{\mathbf{x}}\|$, which we will further assume. In words, k -means tries to find the partitioning γ_k^* that minimizes the sums of squared distances of objects to the centroid of their clusters.

The definition of L_k implies that k -means can only produce "convex" clusters with a tendency towards a spherical shape. In fact, the assignment induced by γ_k^* partitions the data by Voronoi cells of the centroids $\{\bar{\mathbf{x}}_1(\gamma_k^*), \dots, \bar{\mathbf{x}}_k(\gamma_k^*)\}$.

An interesting alternative loss function, which has the same set of minimizers, is

$$\tilde{L}_k(\gamma) := \sum_{j=1}^k n_j^{-1}(\gamma) \sum_{i_1, i_2 \in C_j(\gamma)} \|\mathbf{x}_{i_1} - \mathbf{x}_{i_2}\|^2, \quad (19)$$

that is, γ_k^* minimizes a weighted sum of within-cluster dispersion, i.e., the sum of squared pairwise distances between the feature vectors of its members. Besides its interpretational value, the expression (19) does not involve

³⁷There are some variants that understand clustering in a slightly or substantially different ways than just "choose one γ^* from $\Gamma_{n,k}$ ". For instance, the DBSCAN clustering does not require specifying k in advance and designates some objects as "noise". Other variants, not mentioned in this introductory text, include fuzzy clustering, soft assignment clustering and so on.

centroids; hence it can be evaluated directly from the matrix of pairwise distances between the feature vectors.

The computation of the minimum of L_k (or \tilde{L}_k) is in general a difficult problem of discrete optimization.³⁸ In some special cases, the optimal γ_k^* is not unique and there is a large number of local optima of the objective function³⁹, some of which can be very different from what we would intuitively deem to be a good partitioning. However, there exist simple heuristics that typically find a good local minimum quickly.

One such heuristic is the so-called Lloyd's algorithm, which creates a sequence $\gamma_0, \gamma_1, \dots, \gamma_{\text{stop}}$ of partitionings as follows:

1. Set $t \leftarrow 0$ and create γ_0 randomly.⁴⁰
2. Until a given limit t_{max} on t is reached or the loss function is not improved⁴¹ repeat the following steps:
 - (a) Compute the centroids $\bar{\mathbf{x}}_j(\gamma_t)$ for all clusters $j = 1, \dots, k$.
 - (b) Create γ_{t+1} by assigning objects $i = 1, \dots, n$ to the nearest $\bar{\mathbf{x}}_j(\gamma_t)$.⁴²
 - (c) Set $t \leftarrow t + 1$.

Since the Lloyd's algorithm can get stuck in a local optimum, it is common to run it multiple times (with multiple initial clusterings γ_0), and select the one of the final clusterings γ_{stop} that provides the minimum value of L_k . There exist some more efficient (and more complex) heuristics for minimizing L_k than the basic Lloyd's method such as the Hartigan-Wong method.⁴³

While k -means often provides a reasonable clustering, one must keep in mind that its results depend on the units (scales) of the individual components of the feature vectors. For instance, if one variable has a much larger

³⁸Finding the global minimizer of k -means is NP-hard in general.

³⁹Of course, the notion of a local optimum requires that we define what we mean a "neighbourhood" of a partitioning γ , but the claim that the problem can have many local optima is valid for most natural definitions of the structure of neighbourhoods.

⁴⁰There are many straightforward methods of creating an initial random γ_0 . An efficient advanced method of constructing an initial γ_0 is called k -means++.

⁴¹That is, if $L_k(\gamma_t) \geq L_k(\gamma_{t-1})$, which must occur sooner or later, because there is only a finite number of labellings.

⁴²Formally, this means that $\gamma_{t+1}(i) = j$ if $j \in \operatorname{argmin}_{l=1, \dots, k} \|\mathbf{x}_i - \bar{\mathbf{x}}_l(\gamma_t)\|$ with ties resolved randomly. If the current centroid is one of the centroids closest to \mathbf{x}_i , then we do not change the cluster of i .

⁴³This method is in spirit similar to the Fedorov's exchange algorithm that you encountered in the course on Experimental design.

variance than the others, the clustering will be driven almost entirely by that variable, with the remaining variables having only a negligible effect. As in PCA, it is therefore common to standardize the variables to unit variance before applying k -means.

An interesting modification of k -means is the method called k -medoids. Note that the standard k -means algorithms require that we compute the centroids of clusters, which in turn assumes that each object is characterised by a real-valued vector of features. In some cases, however, we only have an $n \times n$ matrix D of mutual “dissimilarities” of all pairs of objects.

The matrix D (with elements $D_{i,j}$) is enough to define the so-called medoid of a set $C \subseteq \{1, \dots, n\}$. A medoid of C is any object $i(C) \in C$ minimizing the average dissimilarity to all other objects of C , i.e.,

$$\sum_{i \in C} D_{i,i(C)} = \min_{l \in C} \sum_{i \in C} D_{i,l}.$$

The method of k -medoids then tries to choose the clustering γ_k^* that minimizes

$$L_k^D(\gamma) := \sum_{j=1}^k \sum_{i \in C_j(\gamma)} D_{i,i(C_j(\gamma))}.$$

Note that the medoids of the optimal clustering γ_k^* can be viewed as the “most typical representatives” of the individual clusters.⁴⁴ Inspecting these objects often provides a more direct, intuitive understanding of the clusters than inspecting the centroids in k -means, which are typically abstract numerical vectors rather than observed objects. In addition, compared to k -means, the partitioning produced by k -medoids is often more robust to outliers.

Computing optimal clustering with respect to the method of k -medoids is a difficult problem, similarly to k -means. An analogue of the Lloyd’s algorithm for is here a method called “partitioning around medoids”.

Both k -means and k -medoids assume that we already know the appropriate number k of clusters. However, in most applications, k is not known in advance⁴⁵. There are several rules of thumb of selecting the “best” k . The most common is an “elbow” diagram, similar to the one that we used for

⁴⁴For a centroid or a medoid of a cluster we sometimes use the term “cluster prototype”.

⁴⁵As an exercise, try to find some problems, where the number k of clusters is known in advance.

selecting an appropriate number of principal components. In this approach, we plot the piece-wise linear function interpolating the points

$$(1, L_1(\gamma_1^*)), (2, L_2(\gamma_2^*)), \dots, (K, L_K(\gamma_K^*))$$

for some reasonably large K . We will then choose the value of k^* that corresponds to an “elbow” in this diagram. An intuitive explanation of this procedure is that increasing the number of clusters from k^* to $k^* + 1$ has only a small effect on the improvement of the “fit” of the data by the centroids.

Naturally we can also plot the elbow diagram for k -medoids, with L_k^D instead of L_k . Nevertheless, a graphical method typically used to find the appropriate k for k -medoids (and provide additional information about the quality of a given k -medoid clustering) is the so-called silhouette plot. This is a simple heuristic method an interested readers are referred to abundant online materials on this topic.

Different methods for identifying an appropriate number of clusters include the so-called gap statistic or various measures of stability and “internal validity”⁴⁶ of a clustering.

5.2 Model-based Clustering

In this approach, the optimal clustering γ^* is defined by means of an underlying statistical model assumption.

One possible view is that the n objects are characterized by feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ that are realizations of random vectors $\mathbf{X}_1, \dots, \mathbf{X}_n$ which are independent, with a distribution that depends on γ , i.e., γ is considered to be a fixed (and unknown) model parameter. From this model-based perspective, we obtain the optimal clustering as an estimate $\hat{\gamma}$ of the parameter γ , typically by the principle of maximum likelihood.

First, let us analyse the following very simple model: $\mathbf{X}_i \sim N_p(\mu_{\gamma_i}, \sigma^2 I_p)$ for all $i = 1, \dots, n$, where $\gamma \in \Gamma_{n,k}$, $\sigma > 0$ is known and I_p denotes the $p \times p$ identity matrix. Here, the unknown model parameters are $\hat{\gamma}, \hat{\mu}_1, \dots, \hat{\mu}_k$ and the parametric space is $\Gamma_{n,k} \times \mathbb{R}^p \times \dots \times \mathbb{R}^p$. The likelihood function is

$$L(\gamma, \mu_1, \dots, \mu_k | \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n f_{\mu_{\gamma_i}}(\mathbf{x}_i), \quad (20)$$

⁴⁶The silhouette plot is an internal validity measure as it is only based on the data itself. In contrast “external validity” measures use are based on some external information or data.

where

$$f_{\mu}(\mathbf{x}) = \frac{\exp\left(-\frac{1}{2\sigma^2}\|\mathbf{x} - \mu\|^2\right)}{(2\pi\sigma^2)^{p/2}}, \quad \mathbf{x} \in \mathbb{R}^p,$$

is the density of $N_p(\mu, \sigma^2 I_p)$. The log-likelihood of (20) is

$$\begin{aligned} \ln L(\gamma, \mu_1, \dots, \mu_k | \mathbf{x}_1, \dots, \mathbf{x}_n) &= \ln \prod_{i=1}^n f_{\mu_{\gamma_i}}(\mathbf{x}_i) = \\ &= \sum_{i=1}^n \left\{ -\frac{p}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\mathbf{x}_i - \mu_{\gamma_i}\|^2 \right\} = A - B \sum_{i=1}^n \|\mathbf{x}_i - \mu_{\gamma_i}\|^2, \end{aligned}$$

where $B > 0$. Therefore, the maximum likelihood estimate $\hat{\gamma}, \hat{\mu}_1, \dots, \hat{\mu}_k$ minimizes $\sum_{i=1}^n \|\mathbf{x}_i - \mu_{\gamma_i}\|^2$.

For a fixed $\gamma \in \Gamma_{n,k}$ we have

$$\sum_{i=1}^n \|\mathbf{x}_i - \mu_{\gamma_i}\|^2 = \sum_{j=1}^k \sum_{i \in C_j(\gamma)} \|\mathbf{x}_i - \mu_{\gamma_i}\|^2 = \sum_{j=1}^k \sum_{i \in C_j(\gamma)} \|\mathbf{x}_i - \mu_j\|^2.$$

Now, recall that the point that minimizes the sum of the squared Euclidean distances to the points of a finite set C is the center of mass of C . That is, minimizing $\sum_{i \in C_j(\gamma)} \|\mathbf{x}_i - \mu_j\|^2$ with respect to μ_j leads to the solution $\hat{\mu}_j = \bar{\mathbf{x}}_j(\gamma)$, which is the centroid of $C_j(\gamma)$. Consequently, to find the maximum likelihood estimate of $\gamma \in \Gamma_{n,k}$, we need to minimize

$$\sum_{j=1}^k \sum_{i \in C_j(\gamma)} \|\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma)\|^2, \quad (21)$$

which is the problem (18) of k -means. Therefore, the clustering based on the simplest meaningful Gaussian model is exactly the k -means clustering. Incidentally, this clarifies why the k -means method has the tendency to create not only spherical clusters, but spherical clusters of similar extent.

Now, we could gradually relax the assumptions on the covariance matrices of the normal distributions of $\mathbf{X}_1, \dots, \mathbf{X}_n$, e.g., taking σ^2 as unknown, taking the covariance matrices equal to $\sigma_{\gamma_1} I_p, \dots, \sigma_{\gamma_n} I_p$ with unknown and possibly different $\sigma_1, \dots, \sigma_k > 0$, etc. However, here we will directly move to the most general assumption: $\mathbf{X}_i \sim N_p(\mu_{\gamma_i}, \Sigma_{\gamma_i})$, where $\gamma \in \Gamma_{n,k}$, $\mu_1, \dots, \mu_k \in \mathbb{R}^p$, and $\Sigma_1, \dots, \Sigma_k \in \mathcal{S}_{++}^p$,⁴⁷ all unknown.

⁴⁷ \mathcal{S}_{++}^p is the set of all positive-definite $p \times p$ matrices.

Let $\Gamma_{n,k}^R$ denote the set of all “regular” clusterings, i.e., clusterings $\gamma \in \Gamma_{n,k}$ such that the variance-covariance cluster matrices $S_1(\gamma), \dots, S_k(\gamma)$ are non-singular. Let

$$\hat{\gamma}, \hat{\mu}_1, \dots, \hat{\mu}_k, \hat{\Sigma}_1, \dots, \hat{\Sigma}_k$$

denote the model parameters that, on the set $\Gamma_{n,k}^R \times \mathbb{R}^p \times \dots \times \mathbb{R}^p \times \mathcal{S}_{++}^p \times \dots \times \mathcal{S}_{++}^p$, maximize the likelihood function

$$L(\gamma, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k | \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n f_{\mu_{\gamma_i}, \Sigma_{\gamma_i}}(\mathbf{x}_i), \quad (22)$$

where

$$f_{\mu, \Sigma}(\mathbf{x}) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right)}{(2\pi)^{p/2} \sqrt{\det(\Sigma)}}; \quad \mathbf{x} \in \mathbb{R}^p$$

is the density of the distribution $N_p(\mu, \Sigma)$ with a positive definite variance-covariance matrix Σ . In the sequel, we will again derive an optimization problem of a simple form which provides the maximum likelihood estimate $\hat{\gamma}$, without the need to explicitly compute $\hat{\mu}_1, \dots, \hat{\mu}_k, \hat{\Sigma}_1, \dots, \hat{\Sigma}_k$.

Consider the log-likelihood function of (22):

$$\begin{aligned} \ln L(\gamma, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k | \mathbf{x}_1, \dots, \mathbf{x}_n) &= \ln \prod_{i=1}^n f_{\mu_{\gamma_i}, \Sigma_{\gamma_i}}(\mathbf{x}_i) \\ &= \sum_{i=1}^n \left[-\frac{p}{2} \ln(2\pi) - \frac{1}{2} \ln \det(\Sigma_{\gamma_i}) - \frac{1}{2} (\mathbf{x}_i - \mu_{\gamma_i})^T \Sigma_{\gamma_i}^{-1} (\mathbf{x}_i - \mu_{\gamma_i}) \right] \\ &= \sum_{j=1}^k \sum_{i \in C_j(\gamma)} \left[-\frac{p}{2} \ln(2\pi) - \frac{1}{2} \ln \det(\Sigma_j) - \frac{1}{2} (\mathbf{x}_i - \mu_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \mu_j) \right]. \end{aligned}$$

For any fixed $\gamma \in \Gamma_{n,k}^R$ and $j \in \{1, \dots, k\}$ the sum

$$\sum_{i \in C_j(\gamma)} \left[-\frac{p}{2} \ln(2\pi) - \frac{1}{2} \ln \det(\Sigma_j) - \frac{1}{2} (\mathbf{x}_i - \mu_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \mu_j) \right],$$

attains the maximum for $\hat{\mu}_j = \bar{\mathbf{x}}_j(\gamma)$ and $\hat{\Sigma}_j = S_j(\gamma)$ ⁴⁸. Moreover, the constant $-\frac{p}{2} \ln(2\pi)$ plays no role in our optimization problem, therefore $\hat{\gamma}$ is

⁴⁸The proof of this statement is analogous to the proof that the sample mean and the sample variance-covariance matrix are the maximum likelihood estimates based on the random sample from the multivariate normal distribution. You should know this from the course on Classical multivariate statistics led by Assoc. Prof. Lenka Filová

the solution of the problem:

$$\operatorname{argmax}_{\gamma \in \Gamma_{n,k}^R} \sum_{j=1}^k \left[-n_j(\gamma) \ln \det(S_j(\gamma)) - \sum_{i \in C_j(\gamma)} (\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma))^T S_j^{-1}(\gamma) (\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma)) \right] \quad (23)$$

Using the basic properties of the trace⁴⁹ we obtain that for any $\gamma \in \Gamma_{n,k}^R$ and $j \in \{1, \dots, k\}$:

$$\begin{aligned} & \sum_{i \in C_j(\gamma)} (\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma))^T S_j^{-1}(\gamma) (\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma)) \\ &= \sum_{i \in C_j(\gamma)} \operatorname{tr} [(\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma))^T S_j^{-1}(\gamma) (\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma))] \\ &= \operatorname{tr} \left[S_j^{-1}(\gamma) \sum_{i \in C_j(\gamma)} (\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma)) (\mathbf{x}_i - \bar{\mathbf{x}}_j(\gamma))^T \right] \\ &= \operatorname{tr} [S_j^{-1}(\gamma) n_j(\gamma) S_j(\gamma)] = \operatorname{tr} [n_j(\gamma) I_p] = p n_j(\gamma), \end{aligned}$$

which means that the optimization problem (23) is equivalent to

$$\operatorname{argmin}_{\gamma \in \Gamma_{n,k}^R} \sum_{j=1}^k n_j(\gamma) \ln \det(S_j(\gamma)) + \frac{p}{2} \sum_{j=1}^k n_j(\gamma).$$

Clearly, $\sum_{j=1}^k n_j(\gamma) = n$, which does not depend on γ , therefore

$$\hat{\gamma} = \operatorname{argmin}_{\gamma \in \Gamma_{n,k}^R} \sum_{j=1}^k n_j(\gamma) \ln \det(S_j(\gamma)). \quad (24)$$

Note that, roughly speaking, the optimal clustering is the one that minimizes a weighted average of logarithms of “volumes” of clusters⁵⁰. The form of the problem (24) is very natural, simple, but it is again generally very difficult to reliably find its global solution. Here we again need in practice to apply optimization heuristics.

A natural alternative statistical view of how $\mathbf{x}_1, \dots, \mathbf{x}_n$ are generated is as follows: We add a probabilistic mechanism for how the label of each object

⁴⁹We use the fact that a scalar (a 1 by 1 matrix) is trivially the trace of the scalar, then the fact that $\operatorname{tr}(AB) = \operatorname{tr}(BA)$ if both matrix products are compatible and also linearity of the trace. These properties were explained for instance by Dr. Samuel Rosa within the lecture Matrix Algebra for Statisticians.

⁵⁰As an exercise, justify the view that $\det(S_j(\gamma))$ is related to the “volume” of $C_j(\gamma)$.

is generated (i.e., we do not treat the label as a fixed parameter, but as a result of random sampling). More precisely, the labels $\gamma_1, \dots, \gamma_n$ are independent realizations of random variables G_1, \dots, G_n with a common discrete probability distribution on $\{1, \dots, k\}$ determined by probabilities τ_1, \dots, τ_k summing to 1. These probabilities represent the prevalence of individual clusters in the underlying population.

As a result, $\mathbf{x}_1, \dots, \mathbf{x}_n$ are assumed to be realizations of i.i.d. random vectors $\mathbf{X}_1, \dots, \mathbf{X}_n$ from the so-called Gaussian mixture distribution with density

$$\sum_{j=1}^k \tau_j f_{\mu_j, \Sigma_j}(\mathbf{x}); \quad \mathbf{x} \in \mathbb{R}^p,$$

parametrized by $\tau_1, \dots, \tau_k \in \mathbb{P}_k$,⁵¹ $\mu_1, \dots, \mu_k \in \mathbb{R}^p$, and $\Sigma_1, \dots, \Sigma_k \in \mathcal{S}_{++}^p$. These parameters are estimated from the data, again typically via maximum likelihood.

The maximum likelihood approach for a restricted form of the Gaussian mixture model, namely fixing τ_j to $1/k$ (i.e., assuming that the clusters have the same prevalence in the underlying population) and, simultaneously, fixing Σ_j to $\sigma^2 I_p$ (i.e., assuming spherical and equally dispersed shapes of the clusters), again leads to the k -means clustering.

The algorithm used to solve the ML optimization problem for the original (unrestricted) model is usually the so-called EM (expectation-maximization) method. For each object $i = 1, \dots, n$, the resulting label $\hat{\gamma}_i \in \{1, \dots, k\}$ is the one that maximizes the estimated posterior membership probability

$$P(G_i = j | \mathbf{X}_i = \mathbf{x}_i) = \frac{\hat{\tau}_j f_{\hat{\mu}_j, \hat{\Sigma}_j}(\mathbf{x}_i)}{\sum_{\ell=1}^k \hat{\tau}_\ell f_{\hat{\mu}_\ell, \hat{\Sigma}_\ell}(\mathbf{x}_i)}$$

under the fitted model.

Model-based clustering (under the most general assumptions) is more difficult – theoretically as well as computationally – than the standard heuristic clustering methods such as k -means or k -medoids; however, it has several advantages. For instance, the most flexible version of the model-based clustering⁵² is invariant under changes in the scale of variables (changes of the

⁵¹ \mathbb{P}_k denotes the probability simplex in \mathbb{R}^k .

⁵²Note that there are many types of “Gaussian model-based clustering”, in particular variants that impose some restrictions on the underlying variance-covariance matrices of the clusters. Many of these variants are not invariant under variable scaling. Nevertheless, these restricted-covariance models may be needed for data of large dimension p , since the number of free parameters of the unrestricted model grows quadratically with p .

units of measurements), the mixture-of-Gaussians view can give a vector of “degrees of confidence” that an object i belongs to individual clusters, the approach can detect non-circular clusters, “clusters inside other clusters” and so on. The model-based approach also enables computing various information criteria measuring the fit of the model, such as the Akaike information criterion (AIC) or the Bayesian information criterion (BIC), that provide means of choosing a suitable number of clusters.

Besides computational complexity, one disadvantage of the model-based clustering is that it cannot be directly used if we only know the matrix of dissimilarities of the objects.⁵³

5.3 Clustering using DBSCAN

A powerful clustering method of a completely different kind is the so-called density-based spatial clustering of applications with noise (DBSCAN). Here, the clustering is not defined as a solution of an optimization problem, but through a notion of reachability and density-connectedness.

Let $\epsilon > 0$ and $\text{minPts} \geq 3$ be parameters. Let us define the clusters as follows:

1. An object i^* is a “core object” if its ϵ -neighborhood

$$N_\epsilon(\mathbf{x}_{i^*}) := \{i \in \{1, \dots, n\} : \|\mathbf{x}_i - \mathbf{x}_{i^*}\| \leq \epsilon\}$$

contains at least minPts objects (including i^* itself).

2. An object i is “directly reachable” from i^* if i^* is a core object and $\|\mathbf{x}_i - \mathbf{x}_{i^*}\| \leq \epsilon$.
3. An object i_r is “reachable” from an object i_1^* , if there is a sequence $i_1^*, i_2^*, \dots, i_{r-1}^*, i_r$ of objects such that i_{l+1}^* is directly reachable from i_l^* for all $l \in \{1, \dots, r-2\}$, and i_r is directly reachable from i_{r-1}^* .
4. Objects that are not reachable from any core object are classified as “noise”.
5. Two objects i_a and i_b are “density-connected” if there is an object i^* , such that i_a is reachable from i^* and i_b is reachable from i^* .

⁵³Note, however, that we can be creative and suggest the following application of the model-based clustering even in the case that we only know the dissimilarity matrix D : first, use some form of the multidimensional scaling (which uses only D) to *embed* the objects into \mathbb{R}^k . Then apply the model-based clustering to this embedding.

6. It is clear that density-connectedness is a symmetric relation⁵⁴ on the set of “non-noise” objects. Then, a DBSCAN cluster is a maximal density-connected set of non-noise objects.

Let us compare DBSCAN to other clustering methods that we already introduced.

First, an important advantage of DBSCAN is that it can produce a special cluster of “noise” objects, which can be viewed as outliers; this makes DBSCAN robust to outliers. Second, DBSCAN does not require us to choose the number k of clusters; the number of clusters is automatically determined by ϵ and minPts . Third, DBSCAN can identify clusters of various shapes, for instance two interlocked banana-shaped clusters, which is unattainable using the previous methods. Fourth, DBSCAN can be easily modified to work with only a dissimilarity matrix, similarly to k -medoids. Fifth, there exists an efficient algorithm that computes the DBSCAN clustering.⁵⁵

The most important disadvantage of DBSCAN is that it can be very sensitive to the choice of the two parameters; it may be difficult to choose them properly.⁵⁶ Moreover, DBSCAN cannot split the objects into clusters of significantly different densities and the DBSCAN clustering generally depends on the units of variables (units of measurement of individual features); these two problems do not occur with the model-based clustering in the form that we presented in Subsection 5.2. A minor disadvantage is that the result of DBSCAN is not fully deterministic (some “border” objects may end up in different clusters depending on the actual algorithm and the order of the objects in the database).

An algorithm related to DBSCAN is called “OPTICS” (Ordering Points to Identify the Clustering Structure). This algorithm also depends on two parameters, ϵ and minPts ; however, the result of OPTICS is not a split into clusters (plus noise objects) but rather a specific *ordering* of the objects⁵⁷ in the feature space. This ordering is then used to construct the so-called “reachability plot”. We can illustrate the key idea of the algorithm for the case $\epsilon = \infty$ as follows.

⁵⁴It is not, however, an equivalence relation, because it is not transitive. There are some variants of the DBSCAN principle that lead to the relation of equivalence, and therefore a uniquely defined split into the equivalence classes.

⁵⁵We will not describe the algorithm in this text, and for applications it is not in fact necessary to know it (unlike the principle of the method itself).

⁵⁶There is an entire body of literature on choosing ϵ and minPts which goes beyond the scope of the current version of this text.

⁵⁷Again, we will not go into the complex specifics of the algorithm itself.

For $i \in \{1, \dots, n\}$, we define the “core distance” as $d_{core}(\mathbf{x}_i) = \|\mathbf{x}_i - \mathbf{x}_{i^+}\|$, where i^+ is the index of the minPts -th closest point to \mathbf{x}_i . If $i \neq j \in \{1, \dots, n\}$, we define the “reachability distance” $d_{reach}(\mathbf{x}_j | \mathbf{x}_i)$ of \mathbf{x}_j from \mathbf{x}_i as

$$d_{reach}(\mathbf{x}_j | \mathbf{x}_i) = \max\{d_{core}(\mathbf{x}_i), \|\mathbf{x}_j - \mathbf{x}_i\|\}.$$

Briefly formulated, the OPTICS algorithm generates the ordering as follows. The first object is selected at random (or simply as the first object in the database). Then, other objects are added to the ordering in such a way that each newly selected object minimizes the reachability distance from the objects that have already been selected.

Let $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$ be the sequence of feature vectors in the ordering generated by OPTICS. To each $j \in \{2, \dots, n\}$, we assign the reachability distance $r_j = \min_{i < j} d_{reach}(\tilde{\mathbf{x}}_j | \tilde{\mathbf{x}}_i)$. The “reachability plot” is then the plot of the values r_j as a function of $j = 2, \dots, n$. The “valleys” in this plot suggest blocks of points that may form clusters, while the “peaks” indicate transitions between separated clusters.

OPTICS has properties somewhat similar to those of DBSCAN. Compared to DBSCAN, OPTICS is better equipped to identify clusters of varying densities. The reachability plot also provides information on the approximate density of each cluster. However, OPTICS is still not fully invariant under linear transformations of the variables. For instance, the reachability plot (and hence the clustering) may change if we alter the units of measurement for a single variable.

5.4 Using clustering for anomaly detection

Clustering methods can also be used for anomaly detection, especially if the objects can be ordered according to a suitable “degree of anomalousness”. In DBSCAN, the identification of anomalies is built into the method: objects classified as noise can be viewed as anomalous by definition. Moreover, objects that are classified as noise for a wider range of parameter values may be regarded as stronger anomaly candidates. This type of ranking is useful in applications where suspicious cases must be inspected one by one but the inspection capacity is limited, for instance in insurance fraud detection.

Other clustering methods can also provide simple anomaly scores. For example, in k -means or k -medoids, one can order the objects according to the distance from the centroid or medoid of the cluster to which they were assigned; unusually large distances may indicate anomalous objects. Similarly,

very small clusters may correspond to rare or suspicious groups, although they need not always represent anomalies.

In model-based clustering, anomalies can be sought among objects for which the fitted mixture density is very small, or among objects with no clearly dominant posterior membership probability. Thus, clustering methods can often be used not only to identify groups of similar objects, but also to highlight objects that do not fit the main structure of the data well.

Several techniques for outlier and anomaly detection are also based on discrimination and classification methods, which we study next.

6 Classification methods in general

The aim of the classification analysis is to construct a **classifier** \mathcal{C} for **predicting** the **classes** of objects based on the **vectors of features** of the objects.⁵⁸ To this end, we must specify m classes, or categories. Often, the classification is **binary**⁵⁹, that is, $m = 2$, but more than two categories, the so-called **multiclass**⁶⁰ classification, are also common.

The construction of \mathcal{C} is based on the **training data set**, which consists of n objects described by their vectors of features $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ ⁶¹, together with their known correct classifications $y_1, \dots, y_n \in \{1, \dots, m\}$.⁶² are from the The classifications of the objects in the training set are pre-determined by an expert or by another reliable method, which, however, tends to be costly or time consuming. Our goal is to develop a (relatively) inexpensive and rapid classification rule, usually implemented as a computer program, that can categorize future objects using only their known vectors of features⁶³. In

⁵⁸The classes of objects are sometimes called **predicted variables** or **labels**, and the vectors of features are also called **explanatory variables** or **predictors**.

⁵⁹For instance “acceptable/unacceptable”, “female/male”, “infected/noninfected”, etc.

⁶⁰For example the age of a person from m age categories, the soil type from m possible types, a disease from m possible diseases, etc.

⁶¹Note that if we do not directly have real vectors of features of objects, we can usually **embed** the objects into \mathbb{R}^p using various methods. For instance, if we can meaningfully measure the distances between the objects, we can use the multidimensional scaling for the embedding. Also, if the number of features is too large to process, we can use various techniques of **dimensionality reduction** or **feature extraction**.

⁶²Sometimes, the possible labels y_1, \dots, y_n are chosen to be $\{0, \dots, m - 1\}$ or, in the case of a binary classifier, $\{-1, +1\}$.

⁶³Naturally, it is rarely the case that we are able to construct an error-free classification rule, but we can often achieve an acceptably small error rate. The reliability of a classifier is an important topic and will be discussed in the next sections.

the machine learning literature, the process of classification is called **class prediction** or, more generally, **inference**.⁶⁴

Note the distinction between partitioning clustering and classification methods. While clustering methods work only with the data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ representing the vectors of features and try to devise a “natural” split of objects into categories without any additional information⁶⁵, classification methods try to mimic the categorization $y_1, \dots, y_n \in \{1, \dots, m\}$ provided by a “teacher” in order to automatically classify new objects, without further help from the “teacher”. That is why the clustering methods are sometimes called **unsupervised learning**, while the classification methods are called **supervised learning**.

6.1 Representations of classifiers

Often, a classifier is formally represented via real-valued **discrimination functions** l_1, \dots, l_m defined on \mathbb{R}^p . In an abstract way, the value $l_j(\mathbf{x})$ quantifies the “degree of belief” of the classifier that the object with features \mathbf{x} belongs to the class j .⁶⁶ Such a classifier is sometimes called **soft**. If we systematically convert the values $l_1(\mathbf{x}), \dots, l_m(\mathbf{x})$ into a single “winning” category $j \in \{1, \dots, m\}$ that maximizes⁶⁷ $l_j(\mathbf{x})$, we obtain a **hard** classifier.⁶⁸ It is important to add that l_1, \dots, l_m can be very complicated functions, sometimes requiring a complex algorithm to evaluate.

Any hard classifier⁶⁹ can also be represented by a decomposition of \mathbb{R}^p

⁶⁴Note that in the standard mathematical statistics literature, the word “inference” usually means something very different: learning about unknown features of a population or data-generating mechanism from observed data, such as parameter estimation or statistical hypothesis testing.

⁶⁵Except for the *number* of classes. Some clustering techniques, such as the DBSCAN, do not even require the number of classes.

⁶⁶If $l_1(\mathbf{x}), \dots, l_m(\mathbf{x})$ are non-negative and sum to one, they can be viewed as probabilities of individual categories. However, note that the decision that the object be classified into the class i can also be based on other factors, not only on the probability that it originated from the class i , but also on the degree of loss incurred by an incorrect classification.

⁶⁷An observant student could ask what if there are two distinct categories j_1 and j_2 such that both $l_{j_1}(\mathbf{x})$ and $l_{j_2}(\mathbf{x})$ are maximal. In that case we usually select the winning category at random. In this text, we can for simplicity assume that we select the smallest of the maximizing indices. In practice, this does not happen too frequently.

⁶⁸Of course, by converting a soft classifier into a hard classifier we lose some information. The soft classifier is always more informative than the corresponding hard classifier, but in most applications we are ultimately forced to make a definitive decision about the class of the object.

⁶⁹In this study material, we will mostly focus on the hard classifiers.

into m **classification regions** R_1, \dots, R_m , $1 \leq i < j \leq m \Rightarrow R_i \cap R_j = \emptyset$. Here, the object with features \mathbf{x} is classified into the i -th category if and only if \mathbf{x} falls into R_i .

Clearly, any system l_1, \dots, l_m of discrimination functions automatically generates classification regions R_1, \dots, R_m ; a possible choice can be formally defined by $\mathbf{x} \in R_j \Leftrightarrow j = \min(\operatorname{argmax}_i l_i(\mathbf{x}))$, $\mathbf{x} \in \mathbb{R}^p$. The inverse, that is, converting R_1, \dots, R_m into a system of discrimination functions is also simple; we can just use the indicator functions of the regions. However, such a system of discriminant functions does not provide nuanced differences in the degrees of beliefs.

Various classification methods differ in the way they utilize the training set of data to construct either the discrimination functions, or the decomposition R_1, \dots, R_m .

6.2 Reliability of a classifier

The most prominent characteristic of a classifier is its “reliability”.⁷⁰ Assessing the reliability of a classifier is in fact an intricate problem, both theoretically and practically. To put it on a mathematically firm ground, let us adopt the following model of how the data are produced: Let us assume that for each object its class $j \in \{1, \dots, m\}$ is first generated with **prior class probabilities** q_1, \dots, q_m ⁷¹ and then, depending on j (and independent of the previous objects) its vector of features is generated from the **class density** p_j on \mathbb{R}^p .⁷² For this section we will formalize the classifier via the decomposition R_1, \dots, R_m of \mathbb{R}^p .⁷³

Now, the reliability of our classifier can be represented by two related

⁷⁰Many sources use the term “performance” to denote what I call reliability. However, performance is a too general term which would presumably cover for instance the speed of the classification and other characteristics of the classifier. The term reliability indicates that here we only focus on the frequency of errors.

⁷¹Less formally, the probabilities q_1, \dots, q_m represent the frequencies of the m classes.

⁷²It is straightforward to use general probability distributions of features, which could also formally describe discrete (or any) random variables. We choose to only work with densities for simplicity.

⁷³That is, we will only deal with hard classifiers in this section; it more difficult to estimate the precision of a soft classifier. One method of measuring the performance of a predictor which predicts *probabilities* of classes is the so-called Brier score, which is out of the scope of this text.

matrices: First, we could define the matrix with elements

$$P(j|i) = \int_{R_j} p_i(\mathbf{x}) d\mathbf{x}, \quad i, j \in \{1, \dots, m\},$$

which is the probability that an object from class i will be categorized into class j . That is, $P(j|i)$ is the “conditional” probability that the object will be categorized into the class j provided we (only) know that it belongs to the class i .

However, the most commonly used is the so-called **confusion matrix** matrix \mathbf{C} , in which the element $\mathbf{C}_{ij} = P(i, j)$ is the probability that a “generic” object will be from the class i and, at the same time, it will be categorized into the class j .⁷⁴ Clearly, $\mathbf{C}_{ij} = q_i P(j|i)$ for any $i, j \in \{1, \dots, m\}$ and, for a fixed i , the sum of the i th row of \mathbf{C} is the prior frequency of the class i :

$$\sum_{j=1}^m \mathbf{C}_{ij} = \sum_{j=1}^m q_i P(j|i) = q_i \sum_{j=1}^m P(j|i) = q_i.$$

We could also define the column sums of \mathbf{C} , say $r_j = \sum_{i=1}^m \mathbf{C}_{ij}$. The interpretation is that r_j is the probability that a generic object will be *classified* to the class j . The sum of all elements of \mathbf{C} is 1. That is, we can view \mathbf{C} as a (normalized version of a) contingency table.

The off-diagonal elements of \mathbf{C} are called theoretical **miss-classification rates**. Clearly, large values of the diagonal elements of \mathbf{C} and small values of the non-diagonal elements of \mathbf{C} , i.e., small miss-classification rates, signify a reliable classifier. However, there is no universal way of comparing confusion matrices. That is, we often cannot unambiguously decide which of the classifiers A and B is better, based only on their confusion matrices \mathbf{C}_A and \mathbf{C}_B . To this end, we need a real-valued measure defined on the set of possible confusion matrices.

One such measure is based on the expected “loss” from the classification. Let $L(j|i)$ be the estimated loss from classifying an object into the category j , if the correct classification of the object is the category i . Then, the expected loss is $\sum_{i=1}^m \sum_{j=1}^m L(j|i) \mathbf{C}_{ij}$.

⁷⁴This is a “normalized” version of the confusion matrix, which can be interpreted via probabilities. Often, a non-normalized version of the confusion matrix is used representing the frequencies, i.e., $n\mathbf{C}$, where n is the number of objects. In practice, an estimate of \mathbf{C} or $n\mathbf{C}$ is used, as described in one of the following subsections.

If we are unsure about the losses, we can set $L(i|i) = 0$ for all i and $L(j|i) = 1$ for all $i \neq j$. The measure

$$A = 1 - L = \sum_{j=1}^m C_{jj}$$

is called the **accuracy** of the classifier, and can also be interpreted as the probability of the correct classification of a generic new object. Note, however, that sometimes this is a misleading measure.⁷⁵ There are many other possible measures of precision of a classifier, most of which are specialized to binary classification (that is $m = 2$), as discussed next.

6.3 Binary classification

An important special case is the binary classification ($m = 2$). A (hard) binary classifier is fully determined by a decomposition R_1, R_2 , where $R_1 \cup R_2 = \mathbb{R}^p$, $R_1 \cap R_2 = \emptyset$. The boundary between the sets R_1 and R_2 is called the **decision boundary**.

Let us call the two categories “positives” and “negatives”.⁷⁶ It is usual to name the elements of the normalized confusion matrix \mathbf{C} as follows: \mathbf{C}_{11} is the **true positives** rate, \mathbf{C}_{12} is the **false negatives** rate, \mathbf{C}_{21} is the **false positives** rate, and \mathbf{C}_{22} is the **true negatives** rate.⁷⁷ Note that the sum of the true positive and the false negative rates is q_1 , i.e., the overall frequency of the positives in the “population” and the sum of the true negative and the false positive rates is q_2 , i.e., the overall frequency of the negatives in the population.

The ratio of the true positives rate to q_1 , schematically $P(+|+)$ is called the **sensitivity** of the classifier/test. The ratio of the true negatives rate to q_2 , that is, $P(-|-)$, is called the **specificity** of the classifier/test.⁷⁸

⁷⁵For instance if one category is frequent, for instance if $q_1 = 0.9$, then the trivial classifier that classifies all objects into the category 1 has accuracy 0.9.

⁷⁶This is a terminology particularly common in medical applications, where the medical test is in fact a (physical) classifier.

⁷⁷In the Babylon of the machine learning literature, the true positive rate is sometimes called the “probability of detection” and the false positive rate is called the “probability of false alarm”. If instead of the normalized confusion matrix we use a non-normalized version, we simply use the terms true positives, true negatives, false positives, false negatives.

⁷⁸In some fields of application, the authors use the term “recall” or “hit rate” instead of sensitivity and the term “selectivity” instead of specificity.

Besides sensitivity and specificity, there exists a plethora of other measures of the reliability of a classifier. We mentioned some of them, in particular the accuracy, in subsection 6.2. Another popular measure is the **Pearson-Yule phi coefficient**, that is the standard correlation coefficient calculated from \mathbf{C} seen as the probability mass function of two binary random variables.⁷⁹

Let us also mention the so-called “positive predicted value”, or **precision**, formally defined as $\mathbf{C}_{11}/(\mathbf{C}_{11} + \mathbf{C}_{21})$. Precision is analogous to sensitivity: while sensitivity is the probability that an object will be classified to the class 1 provided that it is from the class 1, precision is the probability that an object is from the class 1 provided that it has been classified to the class 1. A good classifier has high values of both sensitivity and precision. Therefore, to quantify the “reliability” of a classifier by a single number, some researchers use the so-called **F1-score**, which is simply the harmonic mean of sensitivity and precision.⁸⁰

For some classification methods, the classifier is fully determined by a single discrimination function l defined on \mathbb{R}^p and a threshold $t \in \mathbb{R}$. That is, the object characterized by \mathbf{x} is classified into the first category if $l(\mathbf{x}) \geq t$ and into the second category if $l(\mathbf{x}) < t$.⁸¹ Compatible classification regions are then $R_1 = \{\mathbf{x} \in \mathbb{R}^p : l(\mathbf{x}) \geq t\}$ and $R_2 = \{\mathbf{x} \in \mathbb{R}^p : l(\mathbf{x}) < t\}$. The decision boundary is then the set $\{\mathbf{x} \in \mathbb{R}^p : l(\mathbf{x}) = t\}$. If l is linear (and the sets R_1 and R_2 are half-spaces and the decision boundary is a hyperplane), such a classifier is usually called **linear**, in the opposite case it is called **non-linear**.

Suppose that the function l is known⁸², but we are allowed to select the threshold t . By changing t we can obtain a continuum of classifiers \mathcal{C}_t , each with its own miss-classification rates. Note that if t decreases, the region R_1 expands. This means that with decreasing t , the rate of the true positives increases (ultimately it reaches 1), but the rate of the false positives also

⁷⁹In machine learning literature, this is called the “Matthews correlation coefficient” in honour of some Matthews, who reinvented the coefficient more than half a century later than Pearson and Yule.

⁸⁰Note that the harmonic mean of two numbers approaches 0 provided that *any* of the two numbers approaches 0.

⁸¹That is, we can formally define the first discriminant function from the previous sections by $l_1(\mathbf{x}) = l(\mathbf{x}) - t$ and the second one by $l_2(\mathbf{x}) = t - l(\mathbf{x})$.

⁸²That is, we assume that the classification method at hand provides the function l per se, not only the partition R_1, R_2 defined by some automatically determined t .

increases (also reaches 1 in the limit).⁸³ Similarly, with increasing t , the rate of the true positives decreases (ultimately it reaches 0), and the rate of the false positives also decreases (reaches 0 in the limit).

Let us now make a plot of all points with coordinates “(1 - specificity, sensitivity)”, formally $(\mathbf{C}_{21}/q_2, \mathbf{C}_{11}/q_1)$.⁸⁴ The set of points will form a curve in the square $[0, 1] \times [0, 1]$ called **receiver operating characteristic curve**⁸⁵, or the **ROC curve**. Note that the ROC curve will be non-decreasing, it starts at $(0, 0)$ and ends at $(1, 1)$.⁸⁶ We can visually inspect the ROC curve and select the most appropriate t , depending on whether we prefer a high sensitivity or a high specificity, or, most likely, a reasonable compromise. Of course, the theoretically ideal situation is if there is some t for which the ROC curve hits the point $(0, 1)$, because it signifies a perfect classifier. While this is usually unattainable in practice, this idea inspires a simple method of measuring the quality of the set of classifiers parametrized by the threshold t (it can be viewed as the quality of the function l itself) by a single number: compute the **area under the ROC curve**, abbreviated by **AUC**. The closer is AUC to 1, the better overall quality of the classifier.⁸⁷

6.4 Estimation of misclassification rates

We could, in principle, calculate the matrix \mathbf{C} from the probabilities q_1, \dots, q_m and the densities p_1, \dots, p_m . However, a fundamental problem is that in practical applications we do not know these characteristics. Therefore, our only hope is to use some method of estimating the elements of \mathbf{C} from the available data. To this end, the most popular methods are the so-called **resampling techniques**.

Suppose that we have a method \mathcal{A} , often a computer algorithm, of constructing a classifier from any training set (\mathbf{X}, \mathbf{y}) .⁸⁸ Here, $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_s)^T$ is the matrix of the feature vectors of the objects in the training set, and $\mathbf{y} = (y_1, \dots, y_s)$ is the vector of “correct” classifications of the objects, sometimes called **labels**.⁸⁹ That is, given any training set (\mathbf{X}, \mathbf{y}) as the input,

⁸³Since q_1 and q_2 are fixed, decreasing t implies increasing the sensitivity but decreasing the specificity.

⁸⁴Note again that $\mathbf{C}_{21}, \mathbf{C}_{11}$ depend on t , although it is not explicit in our notation.

⁸⁵The term is so strange because this curve was first used by electrical engineers.

⁸⁶Except for some pathological cases.

⁸⁷The AUC has also a nice probabilistic interpretation, but we will skip it for now.

⁸⁸ (\mathbf{X}, \mathbf{y}) can be for instance the original, full training set of n objects (in which case $s = n$), or its proper subset.

⁸⁹For the purpose of this section we assume that the labels are $\{1, \dots, m\}$.

the method \mathcal{A} produces some classifier $\mathcal{C} = \mathcal{C}(\mathbf{X}, \mathbf{y}, \mathcal{A})$ represented by a decomposition R_1, \dots, R_m of \mathbb{R}^p . How do we estimate the real misclassification rates of \mathcal{C} applied to future data?

The simplest approach, sometimes called the **in-sample method**, or the “re-substitution” method, is the following: We first apply \mathcal{A} to the entire training dataset to obtain the classifier \mathcal{C} . Next, we go through all objects of the training set one by one and predict their classes using \mathcal{C} ; let us say that the predicted classes are $\tilde{y}_1, \dots, \tilde{y}_n$. Then, the estimate of the (y, \tilde{y}) -entry of \mathbf{C} is the proportion of the pairs (y_i, \tilde{y}_i) , $i = 1, \dots, n$, equal to (y, \tilde{y}) . This method typically yields overly optimistic estimates of \mathbf{C} .⁹⁰

Among the **out-of-sample methods**, the simplest is the **validation-set method**:⁹¹ We first split the full training dataset into two parts, the training part $(\mathbf{X}^t, \mathbf{y}^t)$ and the validation part $(\mathbf{X}^v, \mathbf{y}^v)$, and use \mathcal{A} to create a classifier $\mathcal{C} = \mathcal{C}(\mathbf{X}^t, \mathbf{y}^t, \mathcal{A})$.⁹² Next, we use \mathcal{C} to classify the objects from the validation set.

Let us say that the validation set is composed of the last k objects of the full training set⁹³ and their categories predicted by \mathcal{C} are $\tilde{y}_{n-k+1}, \dots, \tilde{y}_n$. Then, the estimate of the element (y, \tilde{y}) of \mathbf{C} is the proportion of the pairs (y_i, \tilde{y}_i) , $i = n - k + 1, \dots, n$, equal to (y, \tilde{y}) . This method often yields good estimates of the real miss-classification rates of \mathcal{C} . However, its disadvantage is that it is not trained on the full dataset, i.e., it may be less accurate than it could have been if trained on the full dataset.⁹⁴

Finally, we will mention advanced and more computationally demanding out-of-sample estimation methods, called **cross-validation**. In cross-validation, we repeatedly construct classifiers on training subsets derived from the original training dataset and use the corresponding validation subsets to estimate the misclassification rates. Here we rely on systematic aspects of the general method \mathcal{A} that we used for constructing \mathcal{C} .

⁹⁰Intuitively, the reason is that we test the classifier based on the same data that had been used for its training.

⁹¹Some literature uses a different and more nuanced terminology: the validation set is a set of data used in the course of training involving a process called hyper-parameter tuning, and the so-called “test set” is a completely independent dataset used for the evaluation of the final model. Confusingly, sometimes the terms are used in the opposite meaning.

⁹² $(\mathbf{X}^t, \mathbf{y}^t)$ usually contains a fixed proportion, say 80% of objects selected from the full training set.

⁹³Often, the training part is produced by a random selection of $n - k$ objects, but for the notational simplicity we assume here that it is the first $n - k$ objects.

⁹⁴In a sense, we exchanged the precision of the classifier for the precision of the estimation of its miss-classification rates.

A simple cross-validation method is called **leave-one-out**. Here, we repeat the validation-set method n times. In the i -th loop, $i = 1, \dots, n$, we take the object i to be a single-point validation set, construct the classifier \mathcal{C}_i on the training data excluding the object i , and produce an estimate $\hat{\mathbf{C}}_i$ of the matrix \mathbf{C} .⁹⁵ The final estimate of \mathbf{C} is the average of matrices $\hat{\mathbf{C}}_i$. The leave-one-out method generally leads to better estimates of \mathbf{C} than the in-sample method and does not waste training data. However, a big disadvantage is that we need to construct n classifiers. In many applications, the number n of objects is thousands or even millions, and the computation time of this validation method would be prohibitive.

Probably the most frequently used cross-validation method is called the **k -fold cross-validation**, where the validation-set method is repeated k times. We split the training data randomly into k subsets⁹⁶ $(\mathbf{X}^1, \mathbf{y}^1), \dots, (\mathbf{X}^k, \mathbf{y}^k)$. In the t -th loop, $t = 1, \dots, k$, we take $(\mathbf{X}^t, \mathbf{y}^t)$ to be the validation set, construct the classifier \mathcal{C}_t on all training data excluding $(\mathbf{X}^t, \mathbf{y}^t)$, and produce an estimate $\hat{\mathbf{C}}_t$ of the matrix \mathbf{C} . Similarly as in the leave-one-out method, the final estimate of \mathbf{C} is the average of the matrices $\hat{\mathbf{C}}_t$. Clearly, leave-one-out is a special case of the k -fold cross-validation for $k = n$. However, in practice we usually choose k to be much smaller than n .⁹⁷ The k -fold cross-validation seems to be a good compromise between the precision of the estimation of \mathbf{C} and the computation time to produce the estimate.

There are many other variants of cross-validation, for instance leave- p -out, repeated random sub-sampling cross-validation, and so on. Also, some classification methods may have their own, specific, techniques for the estimation of the miss-classification rates.

7 Bayes classifier

7.1 The general Bayes classifier

The practical performance of our classifier can be influenced by the “cost” $L(i|j)$ of miss-classification of an object from the class j into the class i . For example, if $L(i|j)$ is huge⁹⁸ for some pair of classes i, j , then we should be very careful about the miss-classification rate \mathbf{C}_{ji} .

⁹⁵This will be a very coarse estimate with just zeros, except for one number 1 at the position determined by the real and the predicted class of the i -th object.

⁹⁶For simplicity, we assume that n is a multiple of k .

⁹⁷A common choice is $k = 10$.

⁹⁸Suppose that we have a classifier that reads the signal from the traffic lights to inform an autonomous car. The possible classes are now “green”, “yellow” and “red”. Clearly,

For the theoretical development, suppose that we know the underlying probabilistic distributions p_1, \dots, p_m of the feature vectors of objects from categories $1, \dots, m$, and also suppose that we know the probabilities q_1, \dots, q_m of individual classes, that is, q_i is the frequency with which we observe the objects from the class i . In this hypothetical case, it is possible to construct a theoretically optimal classification rule that takes the costs $L(i|j)$ into account, which we will call the **general Bayes classifier**.

To this end, we use the entire matrix $m \times m$ matrix \mathbf{L} of losses from miss-classification, that is $\mathbf{L}_{ii} =: L(i|i) = 0$ for all i and $\mathbf{L}_{ij} =: L(i|j) > 0$ for $i \neq j$, which is the loss that we suffer if we classify an object into the category i , provided that the correct classification of the object is the category j .⁹⁹

For a fixed discrimination rule R_1, \dots, R_m , the probability that an object from category j will be classified into category i is

$$P(i|j) = \int_{R_i} p_j(\mathbf{x}) d\mathbf{x},$$

therefore the probability of the event “the classified object is from the class j and, at the same time, it will be classified to the class i ” is

$$\mathbf{C}_{ji} = q_j \int_{R_i} p_j(\mathbf{x}) d\mathbf{x}.$$

Hence, the mean loss from the classifier R_1, \dots, R_m is

$$E_{R_1, \dots, R_m} = \sum_{i=1}^m \sum_{j=1}^m L(i|j) q_j \int_{R_i} p_j(\mathbf{x}) d\mathbf{x}.$$

Surprisingly, it turns out that it is simple to minimize the mean loss E_{R_1, \dots, R_m} with respect to the choice of the decomposition R_1, \dots, R_m , obtaining thus the theoretical Bayes rule. Indeed, note that

$$E_{R_1, \dots, R_m} = \sum_{i=1}^m \int_{R_i} h_i(\mathbf{x}) d\mathbf{x}, \tag{25}$$

where

$$h_i(\mathbf{x}) = \sum_{j=1}^m L(i|j) q_j p_j(\mathbf{x}) \text{ for all } i = 1, \dots, m$$

the miss-classification “red \rightarrow green” incurs potentially catastrophic costs, therefore the number $L(\text{green}|\text{red})$ is very large.

⁹⁹Note that the matrix \mathbf{L} does not need to be symmetric; it is often the case that the loss $L(i|j)$ is very different from the loss $L(j|i)$ as in the example with the traffic lights.

are fixed functions. Clearly, in order to minimize (25) it is enough to choose any measurable decomposition R_1^*, \dots, R_m^* satisfying

$$R_i^* \subseteq \{\mathbf{x} \in \mathbb{R}^p : h_i(\mathbf{x}) = \min_k h_k(\mathbf{x})\} \text{ for all } i = 1, \dots, m. \quad (26)$$

The selection of optimal rule in (26) can be written more explicitly for two (not mutually exclusive) special cases. First, if $m = 2$, we can choose¹⁰⁰

$$\begin{aligned} R_1^* &= \{\mathbf{x} \in \mathbb{R}^p : h_1(\mathbf{x}) \leq h_2(\mathbf{x})\} \\ &= \{\mathbf{x} \in \mathbb{R}^p : q_2 L(1|2) p_2(\mathbf{x}) \leq q_1 L(2|1) p_1(\mathbf{x})\} \\ &= \left\{ \mathbf{x} \in \mathbb{R}^p : \frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} \geq \frac{q_2 L(1|2)}{q_1 L(2|1)} \right\}, \end{aligned} \quad (27)$$

and $R_2^* = \mathbb{R}^p \setminus R_1^*$. Thus, the binary Bayes classification rule classifies an object with features \mathbf{x} into category 1 if the **likelihood ratio** $p_1(\mathbf{x})/p_2(\mathbf{x})$ exceeds some fixed threshold determined by the prior probabilities q_1, q_2 and losses $L(1|2), L(2|1)$. The likelihood ratio can thus be viewed as the discrimination function l from Section 6.3 and the term $\frac{q_2 L(1|2)}{q_1 L(2|1)}$ as a threshold chosen optimally, provided that the aim is to minimize the expected loss.¹⁰¹

Second, assume a general number m of categories with the simplest losses given by $L(i|j) = 1$ for $i \neq j$.¹⁰² We obtain from (26) that for any $i = 1, \dots, m$:

$$\begin{aligned} R_i^* &\subseteq \{\mathbf{x} \in \mathbb{R}^p : \sum_{j=1, j \neq i}^m q_j p_j(\mathbf{x}) \leq \sum_{j=1, j \neq k}^m q_j p_j(\mathbf{x}) \text{ for all } k\} \\ &= \{\mathbf{x} \in \mathbb{R}^p : q_k p_k(\mathbf{x}) \leq q_i p_i(\mathbf{x}) \text{ for all } k\} \\ &= \{\mathbf{x} \in \mathbb{R}^p : q_i p_i(\mathbf{x}) = \max_k q_k p_k(\mathbf{x})\}. \end{aligned} \quad (28)$$

That is, for a vector \mathbf{x} of features, the Bayes classifier selects the category i such that $q_i p_i(\mathbf{x})$ is maximal. Here, the functions $q_i p_i$ can be viewed as discrimination functions l_i from Section 6.1.

¹⁰⁰We can choose the class for those \mathbf{x} such that $h_1(\mathbf{x}) = h_2(\mathbf{x})$; any such choice leads to the optimal, Bayes, classifier. For simplicity we will also assume that $p_2(\mathbf{x}) > 0$ for all \mathbf{x} .

¹⁰¹However, we can also vary the threshold and obtain the ROC as explained in Section 6.3. The threshold $t = \frac{q_2 L(1|2)}{q_1 L(2|1)}$ will provide a specific point on the ROC.

¹⁰²Note that for these simplest losses the minimization of the expected loss E_{R_1, \dots, R_m} is the same as the maximization of the classifier's accuracy. That is, in this case the Bayes classifier maximizes the accuracy.

For this second special case (i.e., $L(i|j) = 1$ for $i \neq j$), we can also derive the same classifier by the following reasoning. Let Y be the categorical random variable that represents the class of a generic object to be classified in the future. Assume that we observe, for a given object, that the random vector \mathbf{X} representing its features attained the vector $\mathbf{x} \in \mathbb{R}^p$. Then, the posterior probability of the event $[Y = i]$ is given by the Bayes theorem as

$$P(Y = i | \mathbf{X} = \mathbf{x}) = \frac{q_i p_i(\mathbf{x})}{\sum_{k=1}^m q_k p_k(\mathbf{x})}. \quad (29)$$

This means that the (hard) classifier (28) is exactly the one that selects the class with the maximum posterior probability; that is why it is also called the **maximum a posteriori classifier**. The functions at the right-hand side of (29) can be used for the more nuanced information about the class, i.e., what we called the “soft” classifier.

In practice we do not know the densities p_1, \dots, p_m nor the prior probabilities q_1, \dots, q_m , but we can use several reasonable methods of estimating the densities¹⁰³, as well as the prior probabilities.

7.2 Naive Bayes classifier

The so-called **naive Bayes classifier** makes an assumption that all the variables forming vectors of features are independent.¹⁰⁴ In our notation, this means that, for each i the density p_i is a product of one-dimensional marginal densities $p_i^{(1)}, \dots, p_i^{(p)}$:

$$p_i(x_1, \dots, x_p) = \prod_{\ell=1}^p p_i^{(\ell)}(x_\ell), \text{ for all } i = 1, \dots, m. \quad (30)$$

Admittedly, the assumption that all the features behave as independent random variables is unrealistic in all but the simplest cases. However, as a practical tool, the resulting classifier often performs acceptably well and can be our method of choice, if we prefer simplicity over performance.

The simplicity of the naive Bayes classifier is based on the fact that we only need to estimate the marginal densities to obtain an estimate of the

¹⁰³Often, it is enough to estimate the likelihood ratios $p_i(\mathbf{x})/p_j(\mathbf{x})$, or the position of the boundaries of the sets R_i^* .

¹⁰⁴The independence is understood “conditionally for a given class”, that is, provided that we know the class of an object, all its p features, i.e., explanatory random variables, are stochastically independent.

joint density. The marginal densities can be estimated for instance by the univariate normal distribution

$$\widehat{p}_i^{(\ell)}(x_\ell) = f(x_\ell | \widehat{\mu}_{i,\ell}, \widehat{\sigma}_{i,\ell}^2), \quad (31)$$

where $f(x | \mu, \sigma^2)$ denotes the density of the univariate normal distribution with mean μ and variance σ^2 , and the parameter estimates $\widehat{\mu}_{i,\ell}$, $\widehat{\sigma}_{i,\ell}^2$ are computed in a standard way from the training data.¹⁰⁵ This is sometimes called **the Gaussian naive Bayes classifier**. Less coarse estimates of the marginal densities can be obtained via the **kernel density estimation**.

In the literature and in practice, the naive Bayes classifier is typically applied with an implicit assumption that $L(i|j) = 1$ for $i \neq j$; in that case (28) implies that the object characterized by $\mathbf{x} = (x_1, \dots, x_p)$ is classified to the category

$$i^* = \operatorname{argmax}_i \left\{ q_i \prod_{\ell=1}^p \widehat{p}_i^{(\ell)}(x_\ell) \right\}. \quad (32)$$

Note that it is simple to modify the naive Bayes classifier as explained above to the case of discrete (even categorical) features.

7.3 Linear and quadratic discriminant analyses

Another classical simplification is to assume that the densities p_1, \dots, p_m of feature vectors are full-blooded multivariate normal, but all with the same, non-singular variance-covariance matrix Σ . That is,

$$p_i(\mathbf{x}) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \mu_i)^T \Sigma^{-1}(\mathbf{x} - \mu_i)\right)}{(2\pi)^{p/2} \sqrt{\det(\Sigma)}} \quad \text{for all } \mathbf{x} \in \mathbb{R}^p,$$

where $\mu_1, \dots, \mu_m \in \mathbb{R}^p$ are the mean value vectors of distributions p_1, \dots, p_m . Under this assumption, the sets R_1, \dots, R_m turn out to be (possibly unbounded) polyhedrons¹⁰⁶. We will describe the case of two categories in more detail.

¹⁰⁵That is: take the set of the values of the ℓ -th explanatory variable (the ℓ -th feature) of all objects from the i -th class. Compute the standard mean and the standard sample variance of this set of real numbers.

¹⁰⁶In some cases, the optimal choice of R_1, \dots, R_m is not unique, but the sets can always be chosen as polytopes. The situation is, however, markedly different if the variance-covariance matrices of the normal distributions can be different.

Let $m = 2$. The log-likelihood ratio from (27) is for any $\mathbf{x} \in \mathbb{R}^p$

$$\begin{aligned} \ln \left(\frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} \right) &= \frac{1}{2}(\mathbf{x} - \mu_2)^T \Sigma^{-1}(\mathbf{x} - \mu_2) - \frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1) \\ &= (\mu_1 - \mu_2)^T \Sigma^{-1} \mathbf{x} + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1. \end{aligned}$$

Thus, if we denote¹⁰⁷ $\mathbf{a} := \Sigma^{-1}(\mu_1 - \mu_2)$, $b := \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1$, and $k := \frac{q_2 L(1|2)}{q_1 L(2|1)}$, we obtain that

$$R_1^* = \{ \mathbf{x} \in \mathbb{R}^p : \mathbf{a}^T \mathbf{x} + b - \ln(k) \geq 0 \}, \quad R_2^* = \mathbb{R}^p \setminus R_1^*. \quad (33)$$

We obtain a linear classification rule, i.e., R_1^* and R_2^* are half-spaces and the decision boundary is a hyperplane.

Some special cases of the classifier (33) are worth a special attention. First, if the prior probabilities of the two classes are the same, i.e., $q_1 = q_2 = 0.5$ and the losses are symmetric, that is, $L(1|2) = L(2|1)$, then $k = 1$ and $\ln(k) = 0$ in (33). If, moreover, $\Sigma = \sigma^2 I_p$ for some $\sigma^2 > 0$, the hyperplane that sets R_1^* and R_2^* apart passes through the center of mass $(\mu_1 + \mu_2)/2$ of the mean value vectors μ_1, μ_2 , and is orthogonal to $\mu_1 - \mu_2$, i.e., an object is classified into category 1 (or 2) if its feature vector \mathbf{x} is closer to μ_1 than to μ_2 (or closer to μ_2 than to μ_1).

The expression (33) allows us to find simple explicit formulas for the misclassification probabilities $P(2|1)$ and $P(1|2)$. Let $\mathbf{X}_1 \sim N_p(\mu_1, \Sigma)$. Observing that $\mathbf{a}^T \mathbf{X}_1 \sim N(\mathbf{a}^T \mu_1, \mathbf{a}^T \Sigma \mathbf{a})$ we obtain

$$\begin{aligned} P(2|1) &= \int_{R_2} p_1(\mathbf{x}) d\mathbf{x} = P[\mathbf{a}^T \mathbf{X}_1 < \ln k - b] = \\ &= P \left[\frac{\mathbf{a}^T \mathbf{X}_1 - \mathbf{a}^T \mu_1}{\sqrt{\mathbf{a}^T \Sigma \mathbf{a}}} < \frac{\ln k - b - \mathbf{a}^T \mu_1}{\sqrt{\mathbf{a}^T \Sigma \mathbf{a}}} \right] = \Phi \left(\frac{\ln k - b - \mathbf{a}^T \mu_1}{\sqrt{\mathbf{a}^T \Sigma \mathbf{a}}} \right), \end{aligned}$$

where Φ is the distribution function of the standardized normal. Denoting $\alpha = \sqrt{(\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2)}$, which is sometimes called the **Mahalanobis distance** of μ_1 and μ_2 , we can express the misclassification probability as

$$P(2|1) = \Phi \left(\frac{\ln k - \alpha^2/2}{\alpha} \right). \quad (34)$$

An analogous derivation yields

$$P(1|2) = 1 - \Phi \left(\frac{\ln k + \alpha^2/2}{\alpha} \right). \quad (35)$$

¹⁰⁷We assume that $\mu_1 \neq \mu_2$.

Naturally, $P(1|1) = 1 - P(2|1)$ and $P(2|2) = 1 - P(1|2)$ therefore we can obtain all 4 elements of \mathbf{C} from $\mathbf{C}_{ij} = q_i P(j|i)$.

As an interesting exercise, one can show that for the general number m of classes and the special case of $L(i|j) = 1$ for $i \neq j$ and $\Sigma = \sigma^2 I_p$, the classification sets R_i are the **Voronoi cells** of the expected value vectors μ_1, \dots, μ_k . That is, the corresponding classifier selects the category for \mathbf{x} based on the nearest expected value vector.

In practice, the parameters μ_1, \dots, μ_m and Σ are estimated from the training set as follows: $\hat{\mu}_i$ is the mean of the feature vectors assigned to the class y_i in the training set and \mathbf{S} , the estimate of Σ is **pooled** from the matrices $\mathbf{S}_1, \dots, \mathbf{S}_m$ via

$$\mathbf{S} = \frac{(n_1 - 1)\mathbf{S}_1 + \dots + (n_m - 1)\mathbf{S}_m}{n_1 + \dots + n_m - m},$$

where n_i is the number, and \mathbf{S}_i is the sample covariance matrix of the objects assigned to the class y_i in the training set.¹⁰⁸

Note that with the assumption $\Sigma = \sigma^2 I_p$ mentioned above, the finite-sample analogue of theoretical the Voronoi-cell classifier is the **nearest centroid classifier**.

Even if the assumptions of normality are not satisfied¹⁰⁹, the optimal Bayesian approach provides *some* classifier which may perform reasonably well; we can estimate its miss-classification rates using the methods from Section 6.4.¹¹⁰

The methods described above lead to the so-called **Linear discriminant analysis** (LDA), because the discriminant function in (33) is linear in \mathbf{x} . More complex classification methods can be obtained by relaxing the rather stringent assumptions of multivariate normality with the same covariance matrix. For instance, we can assume that the classes have non-unique covariance matrices, i.e., the data are generated from a general mixture of normals. Then we obtain the method called the **Quadratic discriminant**

¹⁰⁸Note once again that here we assume that the covariance matrices of all feature distributions p_i are the same. We also assume that each class has at least 2 representatives in the training set.

¹⁰⁹Which is usually the case in real applications.

¹¹⁰If we have a binary classification problem and we are sure about the theoretical assumptions, we can use $\hat{\mu}_i$ and \mathbf{S} to estimate the Mahalanobis distance of μ_1 and μ_2 , and then use the formulas (35) and (34) to estimate the probabilities of false classification.

analysis (QDA). Here, we separately estimate the means as well as covariance matrices of the explanatory vectors of each class. In the QDA the classification regions R_1, \dots, R_m are in general no longer polytopic, not even in the binary case ($m = 2$).¹¹¹

7.4 K nearest neighbours

K nearest neighbours (KNN) is a simple but powerful approach to constructing a classifier: the object with features \mathbf{x} is classified according to the classes of the k objects of the training set that have features closest to \mathbf{x} . More precisely, for each feature vector \mathbf{x} to be classified, the KNN classifier finds the k closest vectors $\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k}$ from the training set, and classifies \mathbf{x} to the class which occurs most frequently in the set y_{i_1}, \dots, y_{i_k} .¹¹² we call this the **majority vote**. The idea is very simple, but (for a large training set) it may be exceedingly time consuming to find the k nearest neighbours for a given \mathbf{x} .¹¹³

It is interesting to view the KNN principle as a method inspired by the Bayes classification rule. In its simplest multi-class form ($L(i|j) = 1$ for all $i \neq j$ and $q_i = 1/m$ for all i), the Bayes rule (28) says: for \mathbf{x} choose the category j with maximal $p_j(\mathbf{x})$. If the training set is approximately balanced with respect to the classes,¹¹⁴ the KNN is a straightforward method of estimating exactly this - which class has the largest density of features in the point \mathbf{x} .¹¹⁵

There are several other classification methods that can be viewed as approximations to the optimal Bayes classifier. For instance using the kernel estimators of the multivariate densities $p_j, j = 1, \dots, m$, leads to the so-called **kernel density classifier**.

¹¹¹Convince yourself that if $\Sigma_1 \neq \Sigma_2$ then $\ln(p_1(\mathbf{x})/p_2(\mathbf{x}))$ is a quadratic form in \mathbf{x} .

¹¹²If there are several classes that are most frequent in the set y_{i_1}, \dots, y_{i_k} , we select one at random. In the binary classification, we can choose k to be an odd number, to avoid such "ties".

¹¹³There are methods of how to significantly speed up the search for the k -nearest neighbours, especially if we allow occasional errors in their identification. This is a fascinating topic from the point of view of computer science, but we will skip it.

¹¹⁴That is, all classes have approximately same number of objects in the training set, which should be satisfied with the assumption $q_i = 1/m$ for all i , provided that the training set is a representative sample.

¹¹⁵Try to invent alternative methods for predicting which class has the largest density in \mathbf{x} . Use the optimal Bayes classifier to modify the basic KNN to take the probabilities q_i or the miss-classification costs C_{ij} into account.

8 Classification trees

The underlying structure of a standard classification tree¹¹⁶ is that of a directed rooted binary¹¹⁷ tree. There are two main categories of nodes of the classification tree: **non-terminal nodes** and **terminal nodes**. Each non-terminal node is a **parent** of exactly two **descendants**, and each node has exactly one parent, except for the **root node**, which has no parents¹¹⁸. To each non-terminal node we assign a condition (sometimes called a **split**) involving values of a single variable. Moreover, we assign the label “True” to one of the two edges emanating from the non-terminal node, and the label “False” to the other one. To each terminal node we assign a class label from $\{1, \dots, m\}$, where m is the number of classes.¹¹⁹

The process of classification of an object using a classification tree starts in the root node. The object is then passed down the tree via a path directed by the validity (or non-validity) of the conditions in the non-terminal nodes, until the terminal node is hit. The object is assigned to the class corresponding to the label of the terminal node.

Geometrically, the resulting classification R_1, \dots, R_m is formed by p -dimensional polytopic sets (some of them unbounded) with $(p-1)$ -dimensional facets perpendicular to the basic unit vectors $\mathbf{e}_1, \dots, \mathbf{e}_p$.

8.1 Constructing a classification tree using recursive partitioning

The construction of a classification tree is based on the objects in the training set, that is, objects with known true/correct classifications. We will say that a node τ (of a completely or partially constructed tree) “contains” a training set object, if the classification of the training object passes through τ .

The standard procedure for creating a classification tree is called **growing** or **recursive partitioning**¹²⁰ and can be outlined as follows:

¹¹⁶A classification tree is a special case of the structure called a **decision tree**.

¹¹⁷There exist non-binary classification trees but we will omit them in this text.

¹¹⁸Most often, the root node is non-terminal, but for some extreme cases the root node can be terminal, i.e., the classification tree classifies every object into a single category.

¹¹⁹Each terminal node can be assigned a more nuanced output information, for instance frequencies of classes, effectively resulting in the classifier which we call “soft”. Here we will only consider the classification trees that assign just a single class to each terminal node.

¹²⁰Naturally, we could formulate an optimization problem implicitly defining an “optimal”

1. Create a new “empty” node τ which is either the root node or a descendant of an existing non-terminal node.
2. Decide whether τ should be terminal or non-terminal. The node is declared to be terminal if it contains only objects of a single class, or if it contains less than some critical number of training objects, or if the node reached a pre-specified depth (distance from the root node).
3. If τ is terminal, assign a class label to τ . The class label is usually assigned based on the **majority vote**, i.e., the class label is the most frequent label of the objects contained in τ .¹²¹
4. If τ is non-terminal, assign a condition to τ . The common strategy of this step is:
 - (a) Create a finite number of possible conditions that could be assigned to the node τ .
 - (b) Define a “discrimination quality” of each of the conditions for τ (sometimes called the **goodness of split**) and then select the best condition.

Step 4(a) is an easy one, and can be completed before the algorithm is started. For each variable X_i , $i = 1, \dots, p$, do the following:

- If X_i is quantitative or ordinal, select a finite number of thresholds $\theta_i^{(j)}$, $j = 1, \dots, n_i$, and add n_i conditions $X_i < \theta_i^{(1)}, \dots, X_i < \theta_i^{(n_i)}$ to the pool of all permissible conditions.
- If X_i is nominal with possible values from a finite set M_i , add a system of conditions of the form $X_i \in M_i^{(j)}$, $j = 1, \dots, n_i$, where $n_i = 2^{|M_i|-1} - 1$, such that $\emptyset, M_i^{(1)}, \dots, M_i^{(n_i)}, M \setminus M_i^{(1)}, \dots, M \setminus M_i^{(n_i)}, M$ are all $2^{|M|}$ subsets of M .

Thus, we obtained $\sum_{i=1}^p n_i$ permissible conditions to be assigned to τ . The more difficult question is how to construct a natural measure of the discrimination quality of each condition in the node τ .

Let CND be a fixed condition from the pool of all conditions permissible for a non-terminal node τ . Temporarily, assume that CND is assigned to τ .

classification tree; recursive partitioning can then be viewed as a “greedy” method to obtain a possibly sub-optimal but practically often useful classification tree.

¹²¹Of course, the method for assigning the class to the terminal node can also take the expected frequencies of the classes and the miss-classification costs into account.

As the node τ is non-terminal, it will have two descendants, let us denote them as τ_L and τ_R . Suppose that all objects that satisfy *CND* at τ will be directed to τ_L and all objects that do not satisfy *CND* at τ will be directed to τ_R . For the sake of simplicity, we will assume that we only have two classes of objects - *CLS1* and *CLS2* (i.e., $m = 2$). Define the following numbers based on the training set¹²²:

- n_{**} is the total number of objects that pass through τ .
- n_{*1} is the number of objects of *CLS1* that pass through τ .
- n_{*2} is the number of objects of *CLS2* that pass through τ .
- n_{1*} is the number of objects contained in τ that satisfy *CND* (i.e., they will be directed to τ_L).
- n_{2*} is the number of objects contained in τ that do not satisfy *CND* (i.e., they will be directed to τ_R).
- n_{11} is the number of objects of *CLS1* that satisfy *CND*.
- n_{12} is the number of objects of *CLS2* that satisfy *CND*.
- n_{21} is the number of objects of *CLS1* that do not satisfy *CND*.
- n_{22} is the number of objects of *CLS2* that do not satisfy *CND*.

A useful measure of the goodness of the split *CND* at τ can be defined as

$$\iota(\tau) = p_L \iota(\tau_L) + p_R \iota(\tau_R), \quad (36)$$

where $p_L = \frac{n_{1*}}{n_{**}}$, $p_R = \frac{n_{2*}}{n_{**}}$, and $\iota(\tau)$, called an **impurity index** is, loosely speaking, a measure of how unclear the classification would be provided that τ is a terminal node (similarly $\iota(\tau_L)$ and $\iota(\tau_R)$). Two common measures of impurity are the **entropy index** and the so-called **Gini index**. The entropy index is defined as:

$$\iota(\tau) = -\frac{n_{*1}}{n_{**}} \log_2 \frac{n_{*1}}{n_{**}} - \frac{n_{*2}}{n_{**}} \log_2 \frac{n_{*2}}{n_{**}}.$$

Observe that $\iota(\tau)$ is in fact the entropy of the discrete random variable that attains two values with probabilities $\frac{n_{*1}}{n_{**}}$ and $\frac{n_{*2}}{n_{**}}$. If one of the numbers $\frac{n_{*1}}{n_{**}}$ is equal to 1 (and the other number is equal to 0¹²³), then the impurity $\iota(\tau)$

¹²²Draw a contingency table for a clarification.

¹²³By continuity, we define $0 \ln(0) = 0$.

is 0. That is, if all objects in τ are from the same class, it has the minimal possible impurity. On the other hand, if $\frac{n_{*1}}{n_{**}} = \frac{n_{*2}}{n_{**}} = 1/2$ then the impurity $\iota(\tau)$ attains its maximal possible value.

Analogously, we set

$$\begin{aligned}\iota(\tau_L) &= -\frac{n_{11}}{n_{1*}} \log_2 \frac{n_{11}}{n_{1*}} - \frac{n_{12}}{n_{1*}} \log_2 \frac{n_{12}}{n_{1*}}, \\ \iota(\tau_R) &= -\frac{n_{21}}{n_{2*}} \log_2 \frac{n_{21}}{n_{2*}} - \frac{n_{22}}{n_{2*}} \log_2 \frac{n_{22}}{n_{2*}}.\end{aligned}$$

For the Gini index, the impurities are defined by numerically somewhat simpler formulas that approximate the entropy:

$$\begin{aligned}\iota(\tau) &= 1 - \left(\frac{n_{*1}}{n_{**}}\right)^2 - \left(\frac{n_{*2}}{n_{**}}\right)^2, \\ \iota(\tau_L) &= 1 - \left(\frac{n_{11}}{n_{1*}}\right)^2 - \left(\frac{n_{12}}{n_{1*}}\right)^2, \\ \iota(\tau_R) &= 1 - \left(\frac{n_{21}}{n_{2*}}\right)^2 - \left(\frac{n_{22}}{n_{2*}}\right)^2.\end{aligned}$$

Hence, the condition CND^* that maximizes (36) is the one that splits the relevant training set objects into two subsets, such that the classification in these two subsets would be “as clear-cut as possible”.

8.2 Pruning of a classification tree

An efficient strategy of building a classification tree is to grow a “large” tree T_{\max} and then reduce the size of T_{\max} using the process called **pruning**. The basic idea is to define a **pruning measure** of any subtree T of T_{\max} as follows:

$$R_\alpha(T) = R^{err}(T) + \alpha \times \#T,$$

where $R^{err}(T)$ is the re-substitution estimate of the misclassification rate of T , the symbol $\#T$ denotes the number of terminal nodes of T , which serves as a measure of the size of the tree, and $\alpha \geq 0$ is the so-called **complexity parameter**.

For each α , there is one or more “optimal” subtrees of T_{\max} that minimize $R_\alpha(T)$. For $\alpha = 0$, the optimal tree is $T_0 := T_{\max}$ itself and for a very large α , the optimal tree is T_M containing only the root node. By continuously increasing the values of α from 0¹²⁴, we can construct a sequence of optimizing

¹²⁴In a process that requires many technical details that we skip in this text.

trees T_0, T_1, \dots, T_M such that T_i is a subtree of T_{i-1} for all $i = 1, \dots, M$ (such a sequence of trees is called **nested**). From the set T_0, T_1, \dots, T_M we choose the final tree using, for instance, some form of cross-validation.

8.3 Bagging and classification forests

Ensemble methods¹²⁵ combine several individual classifiers into an aggregate classifier. The most important classes of such methods are **bagging**¹²⁶ and **boosting**. In this subsection we will briefly discuss the general method of bagging and a closely related technique of **random forests**, specialized to combine and improve the performance of tree classifiers.

The idea of bagging is as follows. We will construct B ¹²⁷ classifiers $\mathcal{C}_1, \dots, \mathcal{C}_B$, based on randomly selected subsets¹²⁸ $(\mathbf{X}_1, \mathbf{y}_1), \dots, (\mathbf{X}_B, \mathbf{y}_B)$ of the full training dataset (\mathbf{X}, \mathbf{y}) . If \mathbf{x} is the feature vector of a new object to be classified, we classify \mathbf{x} by each of the B classifiers, and base the resulting “ensemble” classification of \mathbf{x} on the majority vote. That is, the classification of the object \mathbf{x} is the class that occurs most frequently in the set of classes predicted by individual classifiers $\mathcal{C}_1, \dots, \mathcal{C}_B$. Bagging improves the performance in particular if the method of construction of the underlying classifiers is “unstable”, i.e., if small perturbations of the training set result in highly variable classifiers. Bagging is typically, but not necessarily, applied to classification trees.

It turns out that the performance of bagging – if applied to classification trees – can be improved by a surprising additional trick: Each split of each bootstrap classification trees $(\mathbf{X}_1, \mathbf{y}_1), \dots, (\mathbf{X}_B, \mathbf{y}_B)$ is based only on a random subset of all explanatory variables (predictors, features, etc). Usually about \sqrt{p} of explanatory variables are randomly selected as candidates for each split. It turns out that this trick makes the individual trees “less correlated”, or more “diverse”, and generally improves the overall performance of the classification. The resulting classifier is called a random forest.

By using bagging or random forests, we lose the simplicity and interpretability of classification trees, but often achieve a classifier with much better performance.

¹²⁵Sometimes a fancy term of “committee machines” is used.

¹²⁶The word is an abbreviation of the expression “bootstrap aggregating”.

¹²⁷ B is typically chosen to be a few thousands.

¹²⁸We use sampling with replacement.

8.4 Boosting

Boosting refers to the methods which combine “weak” classifiers $\mathcal{C}_1, \dots, \mathcal{C}_B$ into one “strong” classifier. These weak classifiers are often, albeit not always, classification trees. In contrast to bagging, the construction is sequential, adding one weak classifier at a time in an optimal way.¹²⁹

Here we will explain the classical version of the method **Adaboost**.¹³⁰ Suppose that our underlying weak classifiers $\mathcal{C}_1, \dots, \mathcal{C}_B$ are binary, and, for any $j = 1, \dots, B$, the classifier \mathcal{C}_j is determined by the discrimination function $k_j : \mathbb{R}^p \rightarrow \{-1, +1\}$.¹³¹ We will also assume that every weak classifier classifies at least one object incorrectly and at least one object correctly.¹³² We will sequentially construct a sequence $\mathcal{S}_1, \mathcal{S}_2, \dots$ of binary classifiers with increasing strength,¹³³ which will be represented by discrimination functions l_1, l_2, \dots .¹³⁴

The initial classifier \mathcal{S}_1 is given by $l_s = k_s$ for some $s \in \{1, \dots, B\}$, i.e., \mathcal{S}_1 is just one of the weak classifiers. Now, suppose that, for a fixed $t \geq 2$, we have already constructed \mathcal{S}_{t-1} determined by its discrimination function l_{t-1} . The main idea of Adaboost is to construct the next classifier \mathcal{S}_t such that its discrimination function is

$$l_t(\mathbf{x}) = l_{t-1}(\mathbf{x}) + \alpha k(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^p,$$

and both $\alpha \in \mathbb{R}$ and $k \in \{k_1, \dots, k_B\}$ are selected to minimize some “error measure” of \mathcal{S}_t . A simple and practically efficient error measure of \mathcal{S}_t is

$$E_t := \sum_{i=1}^n \exp(-y_i l_t(\mathbf{x}_i)). \quad (37)$$

¹²⁹Other differences between boosting and bagging include: 1) the same classifier can be added more than once; 2) the resulting classifier is not based on the majority vote of the weak classifiers, but rather on a special linear combination of the discrimination functions of the weak classifiers. This will be clarified next.

¹³⁰From “adaptive boosting”.

¹³¹That is, we will denote the classes by the labels $-1, +1$. This provides a more convenient notation in mathematical formulas, similarly to the support vector machines described in the next section. Note also that the classifiers $\mathcal{C}_1, \dots, \mathcal{C}_B$ are “hard”; they do not provide nuanced degrees of belief about the class of a classified object.

¹³²Otherwise, the weak classifier would not be weak, or it would be completely useless. Also, these two assumptions lead to a smoother mathematical derivation of Adaboost.

¹³³We will repeat this process a given number of iterations or until some other stopping criterion is satisfied.

¹³⁴In the sense that the sign of $l_t(\mathbf{x})$ determines the class of the object characterized by $\mathbf{x} \in \mathbb{R}^p$, assigned to the object by \mathcal{S}_t .

In (37), y_i is the correct classification of the training object i . Note that if \mathcal{S}_t correctly classifies i , then the contribution of $\exp(-y_i l_t(\mathbf{x}_i))$ to E_t is negligible with $|l_t(\mathbf{x}_i)| \rightarrow \infty$, and if \mathcal{S}_t incorrectly classifies i , the contribution of $\exp(-y_i l_t(\mathbf{x}_i))$ to E_t is enormous if $|l_t(\mathbf{x}_i)| \rightarrow \infty$.

We will express E_t in a form which makes the dependence of the error on α and k more explicit. First, define the weights

$$w_i^{(t)} = \exp(-y_i l_{t-1}(\mathbf{x}_i))$$

for all objects i . The weights are strictly positive numbers, known at the stage of construction of \mathcal{S}_t .¹³⁵ Next, for any $k \in \{k_1, \dots, k_B\}$ let $Q_{yes}^{(k)}$ be the set of all objects i for which $k(\mathbf{x}_i) = y_i$ and let $Q_{no}^{(k)}$ be the set of all objects i for which $k(\mathbf{x}_i) \neq y_i$.¹³⁶ Using this notation we have

$$\begin{aligned} E_t &= \sum_{i=1}^n \exp(-y_i l_{t-1}(\mathbf{x}_i)) \exp(-y_i \alpha k(\mathbf{x}_i)) \\ &= \sum_{i=1}^n w_i^{(t)} \exp(-y_i \alpha k(\mathbf{x}_i)) \\ &= \sum_{i \in Q_{yes}^{(k)}} w_i^{(t)} \exp(-y_i \alpha k(\mathbf{x}_i)) + \sum_{i \in Q_{no}^{(k)}} w_i^{(t)} \exp(-y_i \alpha k(\mathbf{x}_i)). \end{aligned}$$

At this point, we notice that $y_i k(\mathbf{x}_i) = 1$ for all $i \in Q_{yes}^{(k)}$, and $y_i k(\mathbf{x}_i) = -1$ for all $i \in Q_{no}^{(k)}$, hence

$$E_t = \left\{ \sum_{i \in Q_{yes}^{(k)}} w_i^{(t)} \right\} e^{-\alpha} + \left\{ \sum_{i \in Q_{no}^{(k)}} w_i^{(t)} \right\} e^{\alpha}. \quad (38)$$

Now, we need to minimize (38) with respect to the choice of $\alpha \in \mathbb{R}$ and $k \in \{k_1, \dots, k_B\}$. For simplicity, denote the sums of the weights in (38) as

$$w_{yes}^{(k)} := \sum_{i \in Q_{yes}^{(k)}} w_i^{(t)} \quad \text{and} \quad w_{no}^{(k)} := \sum_{i \in Q_{no}^{(k)}} w_i^{(t)}.$$

An important observation is that for any fixed k the function $\alpha \rightarrow w_{yes}^{(k)} e^{-\alpha} + w_{no}^{(k)} e^{\alpha}$ on the right-hand side of (38) is smooth and strictly convex on \mathbb{R} ,

¹³⁵They are known because at this stage the discriminant function l_{t-1} is already known.

¹³⁶Note that for any k both $Q_{yes}^{(k)}$ and $Q_{no}^{(k)}$ are non-empty, which is a consequence of the technical assumption from the introduction of this subsection: each weak classifier classifies at least one object correctly and at least one object incorrectly.

converging to ∞ for $\alpha \rightarrow \infty$ as well as for $\alpha \rightarrow -\infty$;¹³⁷ such a function has a unique minimum on \mathbb{R} given by¹³⁸

$$\alpha^{(k)} = \frac{1}{2} \left(\ln(w_{yes}^{(k)}) - \ln(w_{no}^{(k)}) \right).$$

We can therefore compute all values $E_t = E_t(k, \alpha^{(k)})$, $k = k_1, \dots, k_B$,¹³⁹ and select the smallest one. The arguments of the smallest $E_t(k, \alpha^{(k)})$ provide the optimal weak classifier k^* and the corresponding optimal coefficient α^* .

We remark that typically we have $\alpha^{(k^*)} > 0$, which is equivalent to $w_{yes}^{(k^*)} > w_{no}^{(k^*)}$; that is because the classifiers k , albeit weak, still usually classify the majority of objects correctly. In the case $\alpha^{(k^*)} > 0$ we can characterize k^* more explicitly. To this end, note that (38) with $k = k^*$ and $\alpha = \alpha^{(k^*)}$ can be rewritten to the form

$$E_t = \sum_{i=1}^n w_i^{(t)} e^{-\alpha^{(k^*)}} + w_{no}^{(k^*)} \left(e^{\alpha^{(k^*)}} - e^{-\alpha^{(k^*)}} \right). \quad (39)$$

Now, $\alpha^{(k^*)} > 0$ implies that $e^{\alpha^{(k^*)}} - e^{-\alpha^{(k^*)}} > 0$, therefore k^* must be a weak classifier which provides the minimum value of $w_{no}^{(k)}$.¹⁴⁰ In words, the optimal weak classifier k^* to be incorporated into \mathcal{S}_{t-1} is typically the one which minimizes the sum of the weights $w_i^{(t)}$ of all points i which it classifies incorrectly. Note that the weights $w_1^{(t)}, \dots, w_n^{(t)}$ represent a measure of how “difficult” the objects $i = 1, \dots, n$ are for the classifier \mathcal{S}_{t-1} . Hence, the weak classifier k^* to be added is, roughly put, the one which classifies the most problematic objects correctly, and the better k^* is with respect to this measure of difficulty, the larger is the coefficient $\alpha^{(k^*)}$.

The classifier based on AdaBoost usually performs better than any of the individual weak classifiers, especially when the weak classifiers are only slightly better than the toss of a coin, or when the weak classifiers “complement” each other.

¹³⁷Strict convexity follows from the positivity of the coefficients $w_{yes}^{(k)}$, $w_{no}^{(k)}$ and from the strict convexity of the functions $\alpha \rightarrow e^{-\alpha}$, $\alpha \rightarrow e^{\alpha}$. The limit properties also follow from the positivity of both $w_{yes}^{(k)}$ and $w_{no}^{(k)}$.

¹³⁸It is a simple exercise in calculus.

¹³⁹Let us emphasize that this is only B simple to compute values.

¹⁴⁰We can see this by contradiction. Let there be some k^x , such that $w_{no}^{(k^x)} < w_{no}^{(k^*)}$. Then the pair $k = k^x$, $\alpha = \alpha^{(k^*)}$ would provide a smaller value of E_t than the pair $k = k^*$, $\alpha = \alpha^{(k^*)}$, i.e., the pair $k = k^*$, $\alpha = \alpha^{(k^*)}$ is not the minimizer of E_t .

9 Support vector machines

Support vector machines (SVMs) are binary¹⁴¹ classifiers based on a series of ingenious mathematical ideas as explained in the following subsections.

We will consider a training sample with feature vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ and corresponding classifications $y_1, \dots, y_n \in \{-1, +1\}$,¹⁴² such that none of the feature sets $C_- := \{\mathbf{x}_i : y_i = -1\}$ and $C_+ := \{\mathbf{x}_i : y_i = +1\}$ is empty.

9.1 Linear support vector machines

Assume first that the objects included in the training set are **linearly separable**,¹⁴³ which means that the convex hulls of the feature sets C_- and C_+ are disjoint. An equivalent definition of linear separability is

$$\begin{aligned} y_i = -1 &\Rightarrow \beta_1^T \mathbf{x}_i + \beta_0 < 0, \text{ and} \\ y_i = +1 &\Rightarrow \beta_1^T \mathbf{x}_i + \beta_0 > 0 \end{aligned}$$

for all $i = 1, \dots, n$ and some $\beta = (\beta_0, \beta_1^T)^T \in \mathbb{R}^{p+1}$, where $\mathbf{0}_p \neq \beta_1 \in \mathbb{R}^p$, there exists a hyperplane

$$H_\beta^0 := \{\mathbf{x} \in \mathbb{R}^p : \beta_1^T \mathbf{x} + \beta_0 = 0\}$$

that “separates” C_- and C_+ . Such a **separating hyperplane** is a natural boundary between the classification regions R_{-1} and R_{+1} of the linear classifier that we wish to construct.

Note, however, that in the linearly separable case there is an entire continuum of separating hyperplanes. The fundamental principle of the SVMs is to base the classification on the so-called **maximum-margin** separating hyperplane $H_{\beta^*}^0$. More precisely, define the **margin** of a hyperplane H_β^0 as

$$d_\beta := \min_{i=1, \dots, n} \rho(H_\beta^0, \mathbf{x}_i),$$

¹⁴¹Multi-class SVMs *do* exist, but the fundamental principle of the SVMs pertains to binary classification. One general approach to applying a binary classification method to a problem with $K > 2$ classes is called **one-vs-one**: we construct a collection of $\binom{K}{2}$ binary classifiers between all pairs of classes. A new observation \mathbf{x} is classified by all of the binary classifiers from the collection. The classification of \mathbf{x} is then chosen to be the one that occurs most frequently in the results of these binary classifications.

¹⁴²In the theory of SVMs, the two classes are usually labelled -1 and $+1$ because it permits a more compact statement of the mathematical formulas and results compared to the standard class notations $\{0, 1\}$ or $\{1, 2\}$.

¹⁴³In applications, the objects of the training set are usually *not* linearly separable, but the linearly separable case serves as a springboard to more practical (linearly non-separable and non-linear) versions of the SVMs.

where $\rho(H_\beta^0, \mathbf{x}_i)$ is the Euclidean distance between H_β^0 and \mathbf{x}_i .¹⁴⁴ The optimal $\beta^* = (\beta_0^*, (\beta_1^*)^T)^T \in \mathbb{R}^{p+1}$ is then chosen to maximize d_β subject to the condition that H_β^0 is a separating hyperplane.¹⁴⁵ This fully defines the optimization problem that we need to solve, but the key question is whether we can solve it efficiently enough to be applied to practical problems. We will show that it is indeed the case.

The points \mathbf{x}_i that are closest to $H_{\beta^*}^0$ are called **support vectors** and they play a key role in our optimization problem. Formally, \mathbf{x}_{i^*} is a support vector if and only if

$$\rho(H_{\beta^*}^0, \mathbf{x}_{i^*}) = \min_{i=1, \dots, n} \rho(H_{\beta^*}^0, \mathbf{x}_i) (= d_{\beta^*}).$$

For a fixed β^* , the distance between \mathbf{x} and $H_{\beta^*}^0$ depends only on $|(\beta_1^*)^T \mathbf{x} + \beta_0^*|$, since¹⁴⁶

$$\rho(H_{\beta^*}^0, \mathbf{x}) = \frac{|(\beta_1^*)^T \mathbf{x} + \beta_0^*|}{\|\beta_1^*\|}.$$

Because $H_{\beta^*}^0$ is a maximum-margin separating hyperplane, it resides exactly in the middle between the support vectors from the two opposite classes. Thus, there exists some common $\epsilon > 0$ such that

$$y_i = -1 \Rightarrow (\beta_1^*)^T \mathbf{x}_i + \beta_0^* \leq -\epsilon \quad (40)$$

$$y_i = +1 \Rightarrow (\beta_1^*)^T \mathbf{x}_i + \beta_0^* \geq \epsilon, \quad (41)$$

for all $i = 1, \dots, n$, and the equalities are attained for the support vectors. Consequently,

$$d_{\beta^*} = \min_{i=1, \dots, n} \frac{|(\beta_1^*)^T \mathbf{x}_i + \beta_0^*|}{\|\beta_1^*\|} = \frac{\epsilon}{\|\beta_1^*\|}. \quad (42)$$

Thus, to find the maximum-margin hyperplane, it is enough to find $\beta^* = (\beta_0^*, (\beta_1^*)^T)^T \in \mathbb{R}^{p+1}$ and $\epsilon > 0$ that solve the problem: maximize (42) subject to both (40) and (41).¹⁴⁷ However, if (β^*, ϵ) is a solution to this problem,

¹⁴⁴Margin can be defined in many ways to serve the purpose of introducing the SVMs. Here we define it simply as the distance of H_β^0 to the nearest feature vector.

¹⁴⁵This “geometric” definition already provides some useful insights into the properties of the SVMs. For instance, we see that the SVM classifier can be robust with respect to even extreme outliers. Moreover, the separating hyperplane is unique and well defined, even if $n < p$. These favourable properties are also inherited by the more complex versions of the SVMs that we will study next.

¹⁴⁶You may recall this formula from the class on analytic geometry in the high school.

¹⁴⁷Note that the feasible set for β^* determined by (40) and (41) does not include the zero vector.

then $(\beta^*/\epsilon, 1)$ is also a solution to the same problem,¹⁴⁸ therefore it suffices to consider the problem for the fixed value of ϵ , e.g., $\epsilon = 1$. Finally, note that maximizing $1/\|\beta_1^*\|$ is the same as minimizing $\|\beta_1^*\|^2$ and the pair (40) and (41) of inequalities for $\epsilon = 1$ can be compactly written as $y_i((\beta_1^*)^T \mathbf{x}_i + \beta_0^*) \geq 1$. Therefore, to find a maximum-margin separating hyperplane, it is enough to compute a solution β^* to the following problem of linearly constrained quadratic programming:

$$\begin{aligned} \min_{\beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R}^p} \quad & \|\beta_1\|^2 \\ \text{s.t.} \quad & y_i(\beta_1^T \mathbf{x}_i + \beta_0) \geq 1, \quad i = 1, \dots, n. \end{aligned} \quad (43)$$

This problem is still not computationally trivial, but it is a convex optimization problem¹⁴⁹ which can be solved by robust and efficient solvers of mathematical programming. The linear classifier which we obtain by this method in the case of linearly separable classes is called the **maximum-margin classifier**.

If the class sets C_- and C_+ are **not linearly separable**, the constraints (43) do not admit any feasible solution.¹⁵⁰ In such a case, the idea of how to proceed is to accept violations of the constraints, which is sometimes called the **soft margin** solution. One specification of the soft margin solution is to quantify the violation of the constraints via new non-negative variables, the so-called **slack variables** ξ_1, \dots, ξ_n ; that is, we replace (43) with $y_i(\beta_1^T \mathbf{x}_i + \beta_0) \geq 1 - \xi_i$, $i = 1, \dots, n$. However, the violation should be as small as possible, which can be forced by adding a **penalty term** to the objective function. The standard choice of the penalty term is $C \sum_{i=1}^n \xi_i$,¹⁵¹ where $C > 0$ is a **tuning parameter** chosen by the user of the method. The resulting optimization problem is

$$\begin{aligned} \min_{\beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R}^p, \xi_1 \geq 0, \dots, \xi_n \geq 0} \quad & \|\beta_1\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\beta_1^T \mathbf{x}_i + \beta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n. \end{aligned} \quad (44)$$

¹⁴⁸This is the same as the observation that $H_{\beta^*}^0$, as a set of points, does not change if we divide β^* by the positive ϵ .

¹⁴⁹That is, we minimize a convex function over a non-empty convex set. We remark that the optimization of the parameters of a classifier is a highly non-convex problem for some alternative methods (such as the artificial neural networks).

¹⁵⁰Note that the feasible set (43) is non-empty if and only if the objects of the training set are linearly separable.

¹⁵¹The penalty term should be as simple as possible; with an overcomplicated penalty term, we could render the resulting optimization problem hard to solve.

This optimization problem has always a feasible solution and it is still a problem of linearly constrained quadratic programming, i.e., not that difficult to numerically solve even for relatively large p and n .¹⁵² Now, the optimal β^* will not define a separating hyperplane (there is no such hyperplane in the linearly non-separable case), yet the hyperplane $H_{\beta^*}^0$ forms a reasonable boundary of the classification sets R_{-1} and R_{+1} .

9.2 Nonlinear support vector machines

The practical usability of the linear SVMs is limited to less complex data. Nevertheless, there exists a simple trick of utilizing a linear classification method for constructing a non-linear classifier. The idea is to transform the features via a non-linear transformation¹⁵³ ϕ into a new feature space¹⁵⁴, and construct a linear classifier in that new feature space.¹⁵⁵

For instance, suppose that we replace the original feature vectors \mathbf{x}_i with new feature vectors $\phi(\mathbf{x}_i) := (\mathbf{x}_i^T, \|\mathbf{x}_i\|^2)^T$. Now, if \tilde{R}_-, \tilde{R}_+ is the classification-set representation of a linear classifier¹⁵⁶ in the new feature space \mathbb{R}^{p+1} then the corresponding classifications sets in the original space, i.e., $R_- = \{\mathbf{x} \in \mathbb{R}^p : \phi(\mathbf{x}) \in \tilde{R}_-\}$ and $R_+ = \{\mathbf{x} \in \mathbb{R}^p : \phi(\mathbf{x}) \in \tilde{R}_+\}$, are separated by a non-linear manifold. That is, the corresponding classifier is non-linear.¹⁵⁷

This idea can be used to transform the linear SVMs to non-linear SVMs. Using the **duality theory** of convex optimization¹⁵⁸ it is possible to characterize the solution of (44) via an alternative problem of quadratic optimization as follows.

¹⁵²However, what is usually *really* solved is a dual problem which we formulate in the next subsection.

¹⁵³This transformation is sometimes called a **feature map** in the machine learning literature.

¹⁵⁴For the applications to the SVMs, the new feature space must be a Hilbert space, i.e., a linear vector space with a scalar product $\langle \cdot, \cdot \rangle$. Usually, the dimension of the new feature space is larger than p , often infinite.

¹⁵⁵This trick can be applied to many linear classifiers; however, the SVM classifier is constructed such that it admits an efficient computation, even if the new feature space is huge.

¹⁵⁶That is, \tilde{R}_- and \tilde{R}_+ are complementary half-spaces.

¹⁵⁷Which geometric object corresponds to the boundary between R_- and R_+ ?

¹⁵⁸The theory of duality is challenging to the extent that it is better suited for theoretical lectures on convex optimization. We will thus skip the proof of the following theorem.

Theorem 9.1. Let $C > 0$, let $\beta^* = (\beta_0^*, (\beta_1^*)^T)^T$ be the solution of (44) and let \mathbf{H} be an $n \times n$ matrix with elements $\mathbf{H}_{ij} = y_i y_j \mathbf{x}_i^T \mathbf{x}_j$, $i, j = 1, \dots, n$. Let $\alpha^* \in \mathbf{R}^n$ be the solution to the following optimization problem

$$\begin{aligned} \max_{\alpha \in \mathbf{R}^n} \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \alpha^T \mathbf{H} \alpha \\ \text{s.t.} \quad & \sum_{i=1}^n y_i \alpha_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, n. \end{aligned} \quad (45)$$

Then

$$\begin{aligned} \beta_1^* &= \sum_{i=1}^n y_i \alpha_i^* \mathbf{x}_i, \quad \text{and} \\ \beta_0^* &= y_j - \mathbf{x}_j^T \beta_1^* \quad \text{for any } j \text{ such that } 0 < \alpha_j^* < C. \end{aligned}$$

Note that the matrix \mathbf{H} from Theorem 9.1 is non-negative definite¹⁵⁹ therefore the problem (45) is again one of linearly constrained quadratic programming.¹⁶⁰ However, the really *key* observation is that the problem (45) is based only on the mutual scalar products $\mathbf{x}_i^T \mathbf{x}_j$ between the feature vectors, in the sense that even if we do not know the actual feature vectors, but we do know the mutual scalar products of these vectors, we can solve (45). Specifically, the problem itself does not depend on the dimension p of the features.

To compute the parameter β^* of the hyperplane forming the decision boundary between R_- and R_+ requires that we know the feature vectors \mathbf{x}_i themselves. However, for all practical purposes we do not in fact need to know β^* itself; we only need to be able to classify new feature vectors $\mathbf{x} \in \mathbf{R}^p$ into R_- or R_+ . That is, we need to be able to decide whether the value $f(\mathbf{x}) = (\beta_1^*)^T \mathbf{x} + \beta_0^*$ of the discrimination function corresponding to the separating hyperplane $H_{\beta^*}^0$ is negative or positive. But (for any j such that $0 < \alpha_j < C$) we have

$$f(\mathbf{x}) = (\beta_1^*)^T \mathbf{x} + \beta_0^* = \sum_{i=1}^n y_i \alpha_i^* \mathbf{x}_i^T \mathbf{x} + y_j - \sum_{i=1}^n y_i \alpha_i^* \mathbf{x}_j^T \mathbf{x}_i.$$

Hence, the value of $f(\mathbf{x})$ is *also* fully determined by the scalar products $\mathbf{x}_i^T \mathbf{x}$ and $\mathbf{x}_j^T \mathbf{x}_i$ on the feature space.

¹⁵⁹ \mathbf{H} is a matrix of all scalar products of the set of vectors $y_1 \mathbf{x}_1, \dots, y_n \mathbf{x}_n$, that is, it is the so-called Gram matrix.

¹⁶⁰Note also the strict inequalities in the formula for β_0^* . In fact, the vectors x_j such that $0 < \alpha_j^*$ are called “support vectors” for this “soft margin” solution, although the geometric meaning is less clear than in the linearly separable situation.

Now, let us recall that one method of performing non-linear classification is to first construct an auxiliary linear classifier based on $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)$ for some non-linear feature map ϕ , and then base the classification of \mathbf{x} on how $\phi(\mathbf{x})$ would be classified via this linear classifier. Based on Theorem 9.1 and the explanation above we see that in order to use the linear SVM classifier in the role of this auxiliary linear classifier, we only need to be able to compute the scalar products $\langle \phi(\mathbf{x}), \phi(\tilde{\mathbf{x}}) \rangle$ for any $\mathbf{x}, \tilde{\mathbf{x}} \in \mathbb{R}^p$. We do not need to explicitly compute the vectors $\phi(\mathbf{x})$ at all. We do not even need to know explicitly the new feature space (the codomain of ϕ).

The usual method of non-linear SVM classification is thus to use just a **kernel function**¹⁶¹ $K(\mathbf{x}, \tilde{\mathbf{x}})$ equal to the scalar product $\langle \phi(\mathbf{x}), \phi(\tilde{\mathbf{x}}) \rangle$ of a function ϕ .¹⁶² The feature map ϕ is typically not explicitly stated (nor is the Hilbert space codomain of ϕ).¹⁶³ The standard choices of such kernel functions are the so-called **Gaussian kernel**¹⁶⁴

$$K(\mathbf{x}, \tilde{\mathbf{x}}) = \exp\left(-\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|^2}{2\sigma^2}\right)$$

and the **polynomial kernel**

$$K(\mathbf{x}, \tilde{\mathbf{x}}) = (\mathbf{x}^T \tilde{\mathbf{x}} + c)^d,$$

but there are many other types of kernels. The values $\sigma > 0, c \geq 0, d \in \mathbb{N}$ represent tuning parameters and the software that computes the non-linear SVMs usually allows the user to choose them, similarly to the parameter C .¹⁶⁵

Note that despite the fact that we implicitly work with possibly infinite-dimensional transformed features $\phi(\mathbf{x})$, the size of the quadratic problem from theorem 9.1 is only given by the number n of objects in the training set. However, to construct the problem we need to pre-compute the elements

¹⁶¹Some literature uses the term “support vector machine” only if the classifier uses the kernels. The less advanced versions are then called just a “support vector classifier”.

¹⁶²In fact the features \mathbf{x} do not even need to be real-valued vectors (they can be for instance text strings) provided that we have a method (more specifically, a kernel) that produces scalar products between pairs of objects.

¹⁶³Of course, ϕ and the Hilbert space do exist; we just do not need to explicitly know them for practical purposes.

¹⁶⁴Gaussian kernel is sometimes called the **radial basis kernel**.

¹⁶⁵The choice of these so-called **hyper-parameters** is usually based on the grid-search exploration of the empirical performance of the resulting classifiers. Observe also that if we choose $c = 0$ and $d = 1$ in the polynomial kernel, we have $K(\mathbf{x}, \tilde{\mathbf{x}}) = \mathbf{x}^T \tilde{\mathbf{x}}$. In this case the mapping ϕ is in fact linear - the identity.

$\mathbf{H}_{ij} = y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$ which can take much time if the dimension p is large; therefore, the computational complexity is not completely independent of p . For the evaluation of the discriminant function $f(\mathbf{x})$, the additional computation of the kernel values $K(\mathbf{x}_i, \mathbf{x})$ is required. Nevertheless, the number of non-zero α_i^* 's tends to be much smaller than n in practical problems, and, to classify \mathbf{x} , we need to compute the correspondingly smaller number of values $K(\mathbf{x}_i, \mathbf{x})$.

To summarize, SVMs are classification methods that allow for complex, non-linear separation of classes, yet they can be efficiently computed via methods of quadratic convex optimization.

10 Important properties of classifiers

Although predictive-performance measures of a classifier, such as accuracy, the F1 score, or the expected loss due to misclassification, are very important, several other factors must also be considered when choosing a classifier for a particular application. We focus mainly on the methods introduced in these lecture notes: linear discriminant analysis (LDA), quadratic discriminant analysis (QDA), k -nearest neighbours (kNN), logistic regression (LR), simple classification trees (CT), random forests (RF), AdaBoost with weak trees (AB), support vector machines (SVM), and artificial neural networks (ANN), although ANNs are treated in a separate lecture.

- How difficult is it to construct, prepare, or train the classifier, in terms of hardware resources, time, memory, and energy consumption? Some classifiers require almost no explicit training phase, at least in their simplest form; a typical example is kNN. Other methods are usually computationally cheap to train, such as LDA, QDA, and LR, although difficulties may arise in high dimensions or with ill-conditioned data. Tree-based methods such as CT, RF, and AB require more algorithmic work, especially when many trees are grown. Some classifiers require solving more demanding optimization problems. Typical examples are SVMs, especially nonlinear kernel SVMs, and ANNs, which are usually trained by gradient-based optimization methods and may require very large computational resources for large datasets or complex architectures.
- How difficult is the classification, or inference, of new objects once the classifier is ready to use? A shallow classification tree can be applied very rapidly, even without a computer. Methods such as LDA, LR,

and QDA are also usually cheap to apply. RF and AB are typically still reasonably fast, although their cost increases with the number and depth of the trees. The cost of applying an SVM depends, among other things, on the number of support vectors and on the kernel used. Some methods can be computationally expensive at prediction time, for example a naive kNN classifier with a large training set, or a large deep ANN without suitable hardware.

- How easy or difficult is it to understand the inner workings of the classifier, or its concrete prediction for a given feature vector? Can *we* learn something from the classifier? In some methods, such as LDA, LR, kNN, and simple CTs, we can often directly inspect the logic behind the classifier or behind a particular classification. Other methods, such as QDA, SVMs with nonlinear kernels, RF, AB, and especially large ANNs, may be much more difficult to understand at a human level. These issues are usually discussed under the terms **interpretability** and **explainability**. Classifiers whose decisions are difficult to understand are often called **black-box classifiers**. Such classifiers may be unsuitable in applications where decisions must be justified, audited, or checked for robustness and fairness.
- For which data types and data sizes can the classification method be efficiently applied, or applied at all? Do the predictors need to be represented as real-valued vectors, as in LDA, QDA, LR, SVMs, and ANNs, or can the method naturally handle categorical predictors, as tree-based methods often can?¹⁶⁶ Is the method naturally binary, such as the basic version of the SVM discussed in this lecture, or does it directly handle multiclass classification? Does the method work well with highly imbalanced data, where one class is much less frequent than another, or are class weights, resampling, or other modifications needed?

There are many other considerations when choosing a classifier. Do we need a soft classifier that returns estimated class probabilities, or is a hard classifier sufficient? Can unequal class frequencies, unequal misclassification costs, or object weights be easily incorporated into training? How easy is it to update the classifier when new training data become available? Is the method invariant, or at least insensitive, to changes of units, translations, monotone transformations of predictors, or permutations of features or observations?

¹⁶⁶Some software implementations still require categorical predictors to be encoded numerically.

How sensitive is it to outliers, missing values, noise, and extrapolation outside the range of the training data?

In general, the appropriateness of a classifier depends strongly on the application and on the available training data. It is often impossible to say in advance which classification method will be best. In practice, we usually compare several reasonable methods empirically, using validation procedures that mimic the intended future use of the classifier. As a very rough rule of thumb, simpler methods such as LDA, LR, shallow CTs, linear SVMs, or kNN for small low-dimensional datasets are often preferable for small datasets with simple structure. More complex data, especially large tabular datasets with nonlinearities and interactions, are often better handled by ensemble methods such as RF or AB. Very large and highly structured data, such as images, text, or signals, often require more flexible methods such as ANNs.