

1 Why do we need multivariate statistics?

Multivariate statistical analysis is referring to the analysis of observations which have been observed simultaneously on several variables. This is in contrast to univariate statistical analysis, where we only have observations of a single variable, and we analyze the statistical behavior of these univariate data. In order to analyze multivariate data, we need appropriate tools, called multivariate statistical methods. Multivariate statistics is thus the extension of univariate statistics to more than one dimension

2 What is the Spectral Decomposition Theorem?

Theorem 1.4.3 (Spectral Theorem or eigendecomposition)

Every symmetric matrix Σ of order $(p \times p)$ can be decomposed as

$$\Sigma = \Gamma \mathbf{A} \Gamma^\top = \sum_{i=1}^p a_i \gamma_i \gamma_i^\top, \quad (1.7)$$

where $\mathbf{A} = \text{Diag}(a_1, \dots, a_p)$ is a diagonal matrix with the eigenvalues of Σ , and Γ is an orthogonal matrix (i.e. $\Gamma^\top = \Gamma^{-1}$), where the columns γ_i of Γ , $i = 1, \dots, p$, are standardized eigenvectors of Σ .

3 What is the density of the multivariate normal distribution?

Recall the density of the univariate normal distribution,

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-[(x-\mu)/\sigma]^2/2} \quad -\infty < x < \infty, \quad (1.14)$$

with expectation μ and variance σ^2 .

The density of the multivariate normal distribution is obtained by replacing the distance $(x - \mu)/\sigma$ by a multivariate distance

$$(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) := \text{MD}^2(\mathbf{x}), \quad (1.15)$$

also called (squared) *Mahalanobis distance*. Moreover, the constant $\frac{1}{\sqrt{2\pi}\sigma^2}$ is replaced by a term such that the volume described by the multivariate density function is standardized to 1.

Thus, let $\mathbf{x} = (x_1, \dots, x_p)^\top$ be a p -dimensional random variable with expectation $E(\mathbf{x}) = \boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^\top$ and covariance $\text{Cov}(\mathbf{x}) = \Sigma = [(\sigma_{ij})]$. Assume that Σ is positive definite (denoted by $\Sigma \geq \mathbf{O}$), which is the case if (and only if) all eigenvalues are strictly positive.

4 Name distances for clustering, methods, and their respective objective functions/criteria.

In the following we assume p -variate observations $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$, for $i = 1, \dots, n$, sometimes simply called objects.

A widely used *distance or dissimilarity measure* between the i -th and the j -th object is the **Euclidean distance**

$$d(i, j) = \sqrt{\sum_{k=1}^p (x_{ik} - x_{jk})^2} = \|\mathbf{x}_i - \mathbf{x}_j\|, \quad (2.1)$$

also known under the term L_2 norm distance. An alternative, being less sensitive to outliers, is the **Manhattan distance** (or “city-block” distance), defined as

$$d(i, j) = \sum_{k=1}^p |x_{ik} - x_{jk}| = \|\mathbf{x}_i - \mathbf{x}_j\|_1, \quad (2.2)$$

also called L_1 norm distance.

The distance can be computed for every observation pair, which results in an $n \times n$ distance matrix $\mathbf{D} = [(d_{ij})]$. Clearly, this matrix is symmetric, and the main diagonal elements are all zero.

Hierarchical Clustering Methods:

The “similarity” of the combined pair can be measured, and a “height” is associated with this newly formed class. At the end of the process there is only one single cluster left.

This calls for a new definition of similarity, expressing the distances between a group of observations in one cluster with indexes in the set C_k , and another group in the second cluster with indexes in C_l . The number of observations in each group is n_k and n_l , respectively. Now we can define different distance measures between clusters, which are also reflecting the name of the corresponding clustering algorithms:

- Complete Linkage:

$$\max_{i \in C_k, j \in C_l} d(i, j)$$

The similarity between two clusters is thus defined by that pair of objects from the different clusters which has the biggest distance.

- Single Linkage:

$$\min_{i \in C_k, j \in C_l} d(i, j)$$

Here, the similarity is given by the closest observations from two clusters. Single linkage tends to be unbalanced in the sense that big clusters are quickly combined. This procedure tends to produce many small groups and few large groups. Single linkage is also suitable to detect outliers.

- Average Linkage:

$$\frac{1}{n_k n_l} \sum_{i \in C_k} \sum_{j \in C_l} d(i, j)$$

The similarity is defined as the average of all pairwise distances.

- Centroid method:

Here, one first needs to compute the arithmetic means (vectors) of the observations of each cluster, say $\bar{\mathbf{x}}(C_k)$ and $\bar{\mathbf{x}}(C_l)$. Then the similarity between the two clusters is given by the Euclidean distance of the cluster centers (centroids),

- Ward's method:

The similarity between two clusters is defined as the increase of the variance when merging the two clusters,

$$\frac{\|\bar{\mathbf{x}}(C_k) - \bar{\mathbf{x}}(C_l)\|^2}{1/n_k + 1/n_l}.$$

Those two clusters will be merged where the increase is the smallest possible.

For all of the above methods, the similarity needs to be computed for every pair of clusters, and then the pair with the smallest values of the measure will be merged. The resulting hierarchy can be presented in a **dendrogram**, where the horizontal axis presents the observations and the vertical axis the “height”, which

5 Partitioning Methods

The best known algorithm for creating partitions is the **K-means algorithm**. Consider an index set C_k containing the indexes of the observations of the k -th cluster. One can define the so-called *total point scatter* as

$$T = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n d^2(i, j) = \frac{1}{2} \sum_{k=1}^K \sum_{i \in C_k} \left(\sum_{j \in C_k} d^2(i, j) + \sum_{j \notin C_k} d^2(i, j) \right).$$

It is then possible to decompose T as $T = W(C) + B(C)$ into a *within-cluster* point scatter

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \in C_k} d^2(i, j)$$

and a *between-cluster* point scatter

$$B(C) = \frac{1}{2} \sum_{k=1}^K \sum_{i \in C_k} \sum_{j \notin C_k} d^2(i, j)$$

for a given a cluster partition C . Then, $B(C)$ tends to be large when observations assigned to different clusters are far apart. Thus, it is desirable to find a partition (for fixed K) which maximizes $B(C)$, which is equivalent to minimizing $W(C)$.

For K-means clustering one takes the Euclidean distance $d(i, j) = \|\mathbf{x}_i - \mathbf{x}_j\|$ as the distance measure. Then one can rewrite $W(C)$ from above as

$$W(C) = \sum_{k=1}^K n_k \sum_{i \in C_k} \|\mathbf{x}_i - \bar{\mathbf{x}}_k\|^2, \quad (2.3)$$

where n_k is the number of observations in the k -th cluster, and $\bar{\mathbf{x}}_k$ is the arithmetic mean vector of those observations. Minimizing this criterion means that the n observations are assigned to the K clusters in a way that the (average) distances of the observations from a cluster to their cluster center are minimized. This can be solved by an iterative algorithm:

Step 1. Cluster centers $\mathbf{m}_1, \dots, \mathbf{m}_K$ are initialized. Usually, this is done by randomly selecting K observations as initial cluster centers.

Step 2. Minimize the objective function (2.3) by assigning each observation to the closest cluster center.

Step 3. Minimize (2.3) by replacing the cluster centers from Step 2. by the arithmetic means of the observations per cluster.

Step 4. Repeat Steps 2. and 3. until the assignments do not change.

Each of the Steps 2. and 3. reduce the value of $W(C)$ and thus convergence of this procedure is assured. However, one could end up in a local optimum, and it is thus recommended to re-start the procedure with different random initializations, and take that solution which gives the smallest value $W(C)$.

6 Explain model-based clustering and difficulties that could occur.

Model-based clustering can also be used to obtain a cluster partition, but the result would even give a “probability” for the assignment of an observation to a cluster. As the name indicates, model-based clustering makes use of a statistical model for the shape of the clusters. The standard “model” is multivariate normal distribution, i.e., it is assumed that the cluster has the density of a multivariate normal distribution, with a certain location and covariance. This is a big advantage over K-means clustering, since the cluster shapes can be more flexible. The result of K-means is typically spherically shaped clusters (due to the use of Euclidean distances to the center).

A detailed description of model-based clustering can be found in Fraley and Raftery (2002), and in many other sources of these authors. Assume that the data consist of K clusters, generated by multivariate normal densities with expectation $\boldsymbol{\mu}_k$ and covariance $\boldsymbol{\Sigma}_k$, for $k = 1, \dots, K$. Further, the class probabilities are given by so-called mixing coefficients π_1, \dots, π_K , where $\pi_1 + \dots + \pi_K = 1$. All these parameters are unknown, and they are estimated using the EM (expectation maximization) algorithm. Note that the covariance matrices are $p \times p$ matrices, and for larger p

there are many parameters to estimate from the available data, which can lead to instability. For this reason, the cluster “models” can be simplified, by imposing restrictions on the cluster covariance structures.

The simplest possibility for such restrictions is $\Sigma_k = \sigma^2 \mathbf{I}$, for $k = 1, \dots, K$, where \mathbf{I} is the identity matrix and σ^2 is a parameter for the variance. This would imply that all clusters are spherical, with the same radius. The estimation of the covariances thus reduces to estimating only one parameter, the variance σ^2 . A less restricted covariance structure is $\Sigma_k = \sigma_k^2 \mathbf{I}$, for $k = 1, \dots, K$. In this case, the clusters are still spherical, but their size can be different according to their variance σ_k^2 , which needs to be estimated. Figure 2.3 illustrates different covariance structures.

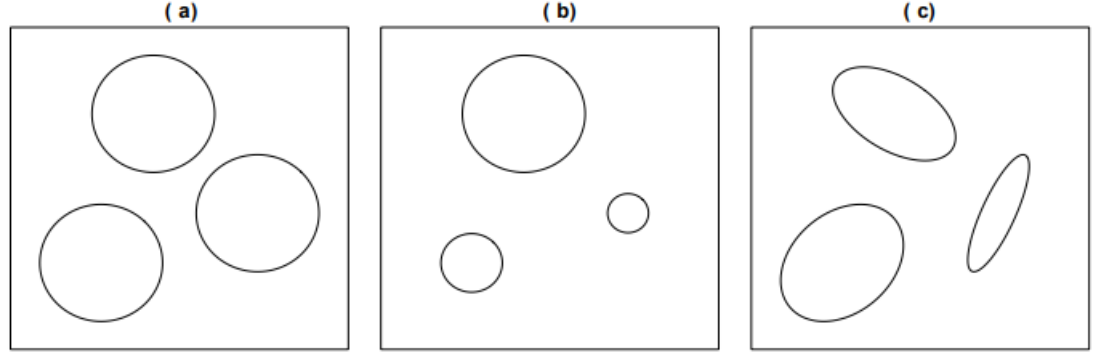


Figure 2.3: Different covariances for three clusters: (a) $\Sigma_1 = \Sigma_2 = \Sigma_3 = \sigma^2 \mathbf{I}$; (b) $\Sigma_j = \sigma_j^2 \mathbf{I}$, for $j = 1, 2, 3$; (c) all Σ_j different and of no special structure.

7 Explain fuzzy clustering, what is the objective function?

Partitioning methods are sometimes called hard clustering methods, since they assign an observation to a cluster or not (1 or 0). In contrast, fuzzy clustering methods.

ods allow for a proportional assignment of an observation to all clusters, where the sum of the proportions is 1. The coefficients for the proportional assignments are called *membership coefficients* $u_{ik} \in [0, 1]$, for $i = 1, \dots, n$ and $k = 1, \dots, K$, with $\sum_{k=1}^K u_{ik} = 1$ for all i .

The best-known fuzzy clustering algorithm is called **fuzzy K-means** algorithm (Bezdek, 1974; Dunn, 1974), and it works very similar to the K-means procedure. First of all, K has to be given. The objective function (2.3) of K-means clustering is replaced by

$$\sum_{i=1}^n \sum_{k=1}^K u_{ik}^2 \|\mathbf{x}_i - \mathbf{m}_k\|^2, \quad (2.4)$$

which has to be minimized. Here, $\mathbf{m}_k = (m_{k1}, \dots, m_{kp})^\top$ is the weighted cluster center of cluster k , defined as

$$m_{kj} = \frac{\sum_{i=1}^n u_{ik}^2 x_{ij}}{\sum_{i=1}^n u_{ik}^2}$$

for $j = 1, \dots, p$.

One can show that the following equality holds:

$$\sum_{i=1}^n \sum_{k=1}^K u_{ik}^2 \|\mathbf{x}_i - \mathbf{m}_k\|^2 = \sum_{k=1}^K \frac{\sum_{i=1}^n \sum_{j=1}^p u_{ik}^2 u_{jk}^2 \|\mathbf{x}_i - \mathbf{x}_j\|^2}{2 \sum_{j=1}^p u_{jk}^2} \quad (2.5)$$

This shows that the squared Euclidean distances between the observations enters the objective function, and thus it is also possible to use the distance matrix rather than the data matrix as an input of the procedure.

8 How can we evaluate clustering solutions -- Hetero/Homogeneity, Calinski-Harabasz, Hartigan, silhouette width, Gap statistic (principles)?

The main goal of cluster analysis is to achieve highly homogeneous clusters, i.e. the observations within a cluster should be very similar to each other. On the other hand, different clusters should be dissimilar, because otherwise they should have been merged into one cluster. In other words, heterogeneity between different clusters should be achieved. Heterogeneity can be measured by

$$B_K = \sum_{k=1}^K \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}\|^2, \quad (2.6)$$

where $\bar{\mathbf{x}}_k$ is the k -th cluster center ($k = 1, \dots, K$), and

$$\bar{\mathbf{x}} = \frac{1}{K} \sum_{k=1}^K \bar{\mathbf{x}}_k$$

is the overall mean of the cluster centers. This term is also called the *between cluster sum of squares*. Homogeneity within the clusters can be defined by

$$W_K = \sum_{k=1}^K \sum_{i \in C_k} \|\mathbf{x}_i - \bar{\mathbf{x}}_k\|^2. \quad (2.7)$$

This term is called the *within cluster sum of squares*, since it considers squared Euclidean distances from the observations to their own cluster center.

While B_K should be large, W_K should be small. However, both measures depend on the number K of clusters, and thus this needs to be considered in a validity measure. Two prominent measures are the **Calinski-Harabasz index**

$$\text{CH}_K = \frac{B_K/(K-1)}{W_K/(n-K)}$$

and the **Hartigan index**

$$\text{H}_K = \ln \frac{B_K}{W_K}.$$

Another prominent validity measure is the **average silhouette width** (Kaufman and Rousseeuw, 1990). Before computing this value, some definitions have to be provided first. The average dissimilarity of an observation \mathbf{x}_i belonging to cluster C_k to all other observations of the same cluster is given by

$$d_{i,C_k} = \frac{1}{n_k - 1} \sum_{i,j \in C_k, i \neq j} d^2(i, j),$$

where n_k is the number of observations in cluster C_k . The average dissimilarity of \mathbf{x}_i to observations from another cluster C_l is given by

$$d_{i,C_l} = \frac{1}{n_l} \sum_{j \in C_l} d^2(i, j).$$

The smallest of these values is

$$d_{i,C} = \min_l d_{i,C_l},$$

and it corresponds to the smallest dissimilarity of the i -th observation to its “closest” cluster. The *silhouette value* is defined as

$$s_i = \frac{d_{i,C} - d_{i,C_k}}{\max(d_{i,C_k}, d_{i,C})}.$$

The values of s_i are within the interval $[-1, 1]$. If the value of s_i is close to 1, the observation is well classified, a value of zero means that the observation is in between two clusters, and a value of -1 refers to a poor classification. Observations with negative silhouette values are probably assigned to a wrong cluster. The average silhouette width is

$$\frac{1}{n} \sum_{i=1}^n s_i,$$

and the higher this value, the better the classification.

Gap Statistic:

The idea is to consider

$$\tilde{W}_K = \sum_{k=1}^K \frac{1}{2n_k} \sum_{i,j \in C_k} d^2(i, j),$$

which is the pooled within-cluster sum of squares around the cluster means. The smaller the value \tilde{W}_K , the more compact are the points in the clusters – but this also depends on K : if K is very big, the clusters are naturally compact. The task is thus to find a possibly small value of K which still yields a small value of \tilde{W}_K .

The Gap statistic is defined as

$$\text{Gap}_n(K) = E_n^*\{\log(\tilde{W}_K)\} - \log(\tilde{W}_K), \quad (2.8)$$

where E_n^* denotes the expectation under a sample of size n from an appropriate reference distribution. For the reference distribution, a unimodal distribution is simulated, generated from a uniform distribution on the hypercube determined by the ranges of the input data. The expectation is estimated by an average of B (e.g. $B = 100$) copies $\log(\tilde{W}_K^*)$, each of which is computed from a bootstrap sample $\mathbf{x}_1^*, \dots, \mathbf{x}_n^*$ of this reference distribution.

Practically, we search for the smallest K such that the Gap statistic yields a local maximum. Since bootstrapping is used, one cannot just compute the average, but also a standard deviation and thus a standard error. Therefore, it is more advisable to look for the smallest K such that the value of the Gap statistic is not more than one standard error away from the first local maximum.

9 What is the least squares estimator?

Based on the observations y_i and x_{i1}, \dots, x_{ip} , for $i = 1, \dots, n$, we want to estimate the regression parameters and the error variance σ^2 by focusing on the differences $y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip}$, which inform about the deviation between the observed response and the model fit. These differences are called **residuals**.

The *least-squares method* minimizes the sum of the squared residuals,

$$\begin{aligned} S(\beta) &= \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 \\ &= (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) . \end{aligned} \quad (3.3)$$

The coefficient minimizing this criterion is denoted as $\hat{\beta}$, and it is the least-squares estimator of the regression parameter β .

We can multiply the terms in Equation (3.3):

$$S(\beta) = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) = \mathbf{y}^\top \mathbf{y} - \beta^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\beta + \beta^\top \mathbf{X}^\top \mathbf{X}\beta . \quad (3.4)$$

The resulting terms are scalars (not vectors or matrices), and thus we have

$$\beta^\top \mathbf{X}^\top \mathbf{y} = (\beta^\top \mathbf{X}^\top \mathbf{y})^\top = \mathbf{y}^\top \mathbf{X}\beta . \quad (3.5)$$

Thus, Equation (3.3) is equal to

$$S(\beta) = \mathbf{y}^\top \mathbf{y} - 2\beta^\top \mathbf{X}^\top \mathbf{y} + \beta^\top \mathbf{X}^\top \mathbf{X}\beta . \quad (3.6)$$

The partial derivative with respect to the vector β yields

$$\frac{\partial S(\beta)}{\partial \beta} = 0 - 2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\beta . \quad (3.7)$$

Setting this equation to zero gives the least-squares estimator

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} . \quad (3.8)$$

Remark: If \mathbf{X} does not have full rank $p + 1 \leq n$, we cannot compute $(\mathbf{X}^\top \mathbf{X})^{-1}$, and a way out would be to use a generalized inverse (e.g. Moore-Penrose) $(\mathbf{X}^\top \mathbf{X})^-$.

With $\hat{\beta}$ one can now compute the *fitted values* $\hat{\mathbf{y}}$ of \mathbf{y} by

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{H}\mathbf{y} , \quad (3.9)$$

where $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is called “hat”-matrix. The *estimated residuals* are

$$\hat{\varepsilon} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{H})\mathbf{y} . \quad (3.10)$$

We can find that $\mathbf{X}^\top \hat{\varepsilon} = \mathbf{X}^\top (\mathbf{I} - \mathbf{H})\mathbf{y} = \mathbf{0}$ and $\hat{\mathbf{y}}^\top \hat{\varepsilon} = \hat{\beta}^\top \mathbf{X}^\top \hat{\varepsilon} = 0$. This means that the estimated residuals are orthogonal to the columns of \mathbf{X} and to the fitted values $\hat{\mathbf{y}}$.

10 How does multivariate linear regression work, what is the objective function, solution and appropriate inference (estimation of covariance of errors). Basic model selection?

While in the multiple linear regression model we considered one response variable, we have m response variables y_1, \dots, y_m in the multivariate regression case. Thus, it is possible to construct for each single response variable a regression model based on the explanatory variables x_1, \dots, x_p :

$$\begin{aligned} y_1 &= \beta_{01} + \beta_{11}x_1 + \dots + \beta_{p1}x_p + \varepsilon_1 \\ \vdots & \quad \quad \quad \vdots \\ y_j &= \beta_{0j} + \beta_{1j}x_1 + \dots + \beta_{pj}x_p + \varepsilon_j \\ \vdots & \quad \quad \quad \vdots \\ y_m &= \beta_{0m} + \beta_{1m}x_1 + \dots + \beta_{pm}x_p + \varepsilon_m \end{aligned} \tag{3.12}$$

For the error term $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_m)^\top$ we now have the assumptions $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ and $\text{Cov}(\boldsymbol{\varepsilon}) = \boldsymbol{\Sigma}$. This means that the error terms for the different responses can be correlated with each other.

Consider now a sample of size n , then we can write down the regression model separately for every response variable $\mathbf{y}_j = (y_{1j}, \dots, y_{nj})^\top$ ($j = 1 \dots, m$) in the same way as in the previous section:

$$\begin{aligned} y_{1j} &= \beta_{0j} + \beta_{1j}x_{11} + \dots + \beta_{pj}x_{1p} + \varepsilon_{1j} \\ y_{2j} &= \beta_{0j} + \beta_{1j}x_{21} + \dots + \beta_{pj}x_{2p} + \varepsilon_{2j} \\ \vdots & \quad \quad \quad \vdots \\ y_{nj} &= \beta_{0j} + \beta_{1j}x_{n1} + \dots + \beta_{pj}x_{np} + \varepsilon_{nj} \end{aligned}$$

All these equations can be collected in matrix form as

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E} \tag{3.13}$$

with the matrices

- $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_m)$ of dimension $n \times m$,
- \mathbf{X} , still of dimension $n \times (p + 1)$,
- $\mathbf{B} = (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m)$ of dimension $(p + 1) \times m$,
- $\mathbf{E} = (\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_m)$ of dimension $n \times m$.

According to the assumptions in the multiple linear regression model, we have

$$E(\varepsilon_j) = \mathbf{0}$$

$$\text{Cov}(\varepsilon_j, \varepsilon_k) = \sigma_{jk} \mathbf{I}_n \quad (3.14)$$

for $j, k = 1, \dots, m$. This means that the m characteristics of the i -th trial ($i = 1, \dots, n$) have covariance $\Sigma = [(\sigma_{jk})]$, but observations from different trials are uncorrelated.

Also in multiple linear regression we are interested here in deriving the least-squares estimator. Similar as in Equation (3.3), we minimize the sum of squared residuals:

$$S(\mathbf{B}) = \sum_{i=1}^n \sum_{j=1}^m (y_{ij} - \beta_{0j} - \beta_{1j}x_{i1} - \dots - \beta_{pj}x_{ip})^2$$

One can show that this is the same as

$$S(\mathbf{B}) = \text{tr}(\mathbf{Y}^\top \mathbf{Y}) - 2 \text{tr}(\mathbf{Y}^\top \mathbf{X} \mathbf{B}) + \text{tr}(\mathbf{B}^\top \mathbf{X}^\top \mathbf{X} \mathbf{B}),$$

where “tr()” denotes the trace of the matrix, i.e. the sum of the diagonal elements. The partial derivative with respect to \mathbf{B} yields

$$\frac{\partial S(\mathbf{B})}{\partial \mathbf{B}} = -2\mathbf{X}^\top \mathbf{Y} + 2\mathbf{X}^\top \mathbf{X} \mathbf{B},$$

and this leads to the least-squares estimator

$$\hat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}. \quad (3.15)$$

Interestingly, this is just the same solution as if we would use the least-squares estimator for each single response variable: Denote $\hat{\mathbf{B}} = (\hat{\beta}_1, \dots, \hat{\beta}_m)$, then we have

$$\hat{\beta}_j = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}_j, \quad (3.16)$$

as the least-squares estimator of β_j for the response \mathbf{y}_j , for $j = 1, \dots, m$.

The fitted values of the responses are obtained by

$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\mathbf{B}}. \quad (3.17)$$

Since the structure of the estimator is still the same as in the multiple linear regression case, still all the results are valid, in particular those of the Gauss-Markov theorem.

11 What are problems with non-robustness? How is it connected to the (empirical) influence function, maxbias curve, breakdown point and efficiency?

The goal of robust statistics is to have less dependency on strict model assumptions. One still works with models and assumptions, but certain deviations and violations are tolerated.

robust methods focus on fitting only the majority of the data, where the requirements need to be fulfilled, but allow for deviations from the minority. Further, covariance estimation is so crucial in multivariate statistics, so that various methods can be robustified just by plugging in a robust covariance estimate.

Good robust estimators downweight outliers during the estimation, but then provide (robust) diagnostics in order to reveal those observations.

Which properties should a robust estimator have? It should be resistant to a sizeable proportion of outliers or deviation from assumptions. It should also still yield reasonable results if these ideal assumptions are valid. In general,

we are interested in the influence function, the maxbias curve and breakdown point, and in the statistical efficiency.

Influence function:

One of the basic ideas of robustness is that a limited amount of contamination should only have a small effect on the estimator. This can be simply empirically checked by varying single data points and looking at the effect on the estimator. We are interested in the effect of one observation when it is varied along the indicated vertical line. For each position of the data point, we are interested in the change of the slope parameter. The result is known as the empirical influence function (EIF). Ideally the EIF should be smooth: it should not show local spikes or should not be a step function:

The concept of the EIF can be formalized by the so-called influence function (IF). The IF measures the influence an infinitesimal amount of contamination has on an estimator with respect to its position in space. More precisely, the influence function of an estimator T at a given distribution G is defined as:

$$\text{IF}(\mathbf{x}, T, G) = \lim_{\varepsilon \downarrow 0} \frac{T[(1 - \varepsilon)G + \varepsilon\delta_{\mathbf{x}}] - T(G)}{\varepsilon}, \quad (4.1)$$

where ε is the fraction of contamination and $\delta_{\mathbf{x}}$ is a probability measure which puts all the mass at \mathbf{x} . For robust estimators, the effect of small contamination will be limited.

Maxbias Curve:

What one expects is that a robust estimator can withstand a certain fraction of contamination. The mathematical tool to examine to which extent an estimator is distorted with respect to the fraction of contamination in the data is the maxbias curve. The maxbias curve measures the bias an estimator has with respect to the percentage of the worst possible type of contamination:

contamination. Let \mathbf{X} be the original data set and $\tilde{\mathbf{X}}$ be a data set in which m out of n observations have been replaced with arbitrary values, and let $\|\cdot\|$ denote the Euclidean norm. Then the maxbias curve for an estimator T is defined as:

$$\text{maxbias}(m, T, \mathbf{X}) = \sup_{\tilde{\mathbf{X}}} \|T(\tilde{\mathbf{X}}) - T(\mathbf{X})\|. \quad (4.2)$$

It is known that for some estimates of regression the worst possible type of outliers is found at points where y , x and the fraction y/x increase to infinity. Non-robust estimators, such as the least-squares regression estimator, turn out to reach a maxbias of infinity already at small amounts of contamination. This brings us to a further important concept.

Breakdown point:

Breakdown point

For every estimators, there exists a point where the maxbias tends to infinity. This point is referred to as the *breakdown point*. Loosely, the breakdown point indicates which percentage of the data may be replaced with outliers before the estimator yields aberrant results. Based on the maxbias curve, the breakdown point of an estimator T at a given sample \mathbf{X} is given by:

$$\varepsilon_n^*(T, \mathbf{X}) = \min \left\{ \frac{m}{n}; \text{maxbias}(m, T, \mathbf{X}) = \infty \right\}. \quad (4.3)$$

For $n \rightarrow \infty$ one obtains the *asymptotic breakdown point*, denoted ε^* . For least-squares regression it holds that $\varepsilon^* = 0$. The maximal possible value of the asymptotic breakdown point of regression estimators equals 0.5 (if one asks for certain equivariance properties).

One of the goals in designing robust estimators is obtaining a high breakdown point. However, bounded influence and high breakdown should not result in a drastic decrease in efficiency.

Statistical efficiency:

The goal is to design robust estimators which are also efficient under normality, but at the same time achieve high (or at least positive) breakdown point and reasonable behavior of the influence function.

12 What are M-estimators? What are the M-estimating equations? Why is it a weighted least squares estimator?

classical least-squares (LS) estimator is defined as

$$\hat{\beta}_{LS} = \arg \min_{\beta} \sum_{i=1}^n r_i(\beta)^2.$$

An ancient alternative to LS is the L_1 estimator defined as

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n |r_i(\beta)|.$$

4.3.1 M-estimators

One can rewrite Equations (4.5) and (4.6) as

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n \rho(r_i(\beta)), \quad (4.7)$$

where $\rho(r) = r^2$ for LS and $\rho(r) = |r|$ for L_1 . By taking other ρ functions, different estimators are obtained. In fact, Equation (4.7) is the definition of a whole class of estimators commonly referred to as *M-estimators*.

In order to not depend on the scale of the response, a more suitable definition of M-estimators of regression is

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n \rho\left(\frac{r_i(\beta)}{\hat{\sigma}}\right), \quad (4.8)$$

where $\hat{\sigma}$ is a robust scale estimator of the residuals, that can be estimated either previously or simultaneously with the regression parameters.

Differentiating (4.8) with respect to β we get that the estimate fulfils the system of *M-estimating equations*

$$\sum_{i=1}^n \psi\left(\frac{r_i(\beta)}{\hat{\sigma}}\right) \mathbf{x}_i = 0 \quad (4.9)$$

where $\psi = \rho'$.

For LS, $\psi(r) = r$, and (4.9) are the well-known normal equations. We may then in general interpret (4.9) as a robustified version of the normal equations, where the residuals are curbed.

Put $W(r) = \psi(r)/r$ and $w_i = W(r_i(\beta)/\hat{\sigma})$. Then (4.9) may be rewritten as

$$\sum_{i=1}^n w_i (y_i - \mathbf{x}_i^\top \beta) \mathbf{x}_i = 0. \quad (4.10)$$

Then (4.10) is a weighted version of the normal equations, and hence the estimator can be seen as weighted LS, with the weights depending on the data. For LS, W is constant. For an estimator to be robust, observations with large residuals should receive a small weight, which implies that $W(r)$ has to decrease to zero fast enough for large r .

The second row of Figure 4.2 refers to the *Huber* family, with ρ' given by

$$\psi(r) = \begin{cases} r & \text{for } |r| \leq b \\ b \operatorname{sign}(r) & \text{otherwise} \end{cases} \quad (4.11)$$

The extreme cases $b \rightarrow \infty$ and $b \rightarrow 0$ correspond to LS and L_1 , respectively.

The last row in Figure 4.2 refers to the Tukey *bisquare* family, with

$$\rho(r) = \begin{cases} \left(\frac{r}{k}\right)^2 \left(3 - 3\left(\frac{r}{k}\right)^2 + \left(\frac{r}{k}\right)^4\right) & \text{for } |r| \leq k \\ 1 & \text{else} \end{cases}. \quad (4.12)$$

When $k \rightarrow \infty$, the corresponding estimate tends to LS and hence becomes more efficient and at the same time less robust. Thus, k is a tuning parameter the choice of which is a compromise between efficiency and robustness. The usual practice is to choose k to attain a given efficiency, such as 0.90.

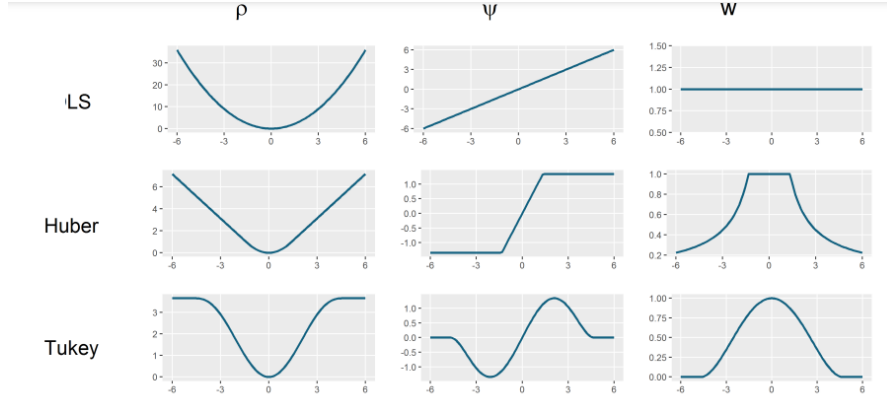


Figure 4.2: Different options for the function ρ , ψ and the weights w .

Note that the Tukey bisquare function has a bounded ρ -function. This is desirable because otherwise if some \mathbf{x}_i is large, then the i -th term will dominate the sum in (4.10), which would be unfortunate if (\mathbf{x}_i, y_i) is atypical (a so-called bad leverage point). For this reason it is better to use M-estimators given by (4.8) with a *bounded* ρ .

13 What are S-estimators? What is the MM-estimator and the properties inherited from S- and M-estimators?

4.3.2 Regression estimators based on a robust residual scale

Given β , let $\mathbf{r}(\beta) = (r_1(\beta), \dots, r_n(\beta))$. We consider an estimator of the form

$$\hat{\beta} = \arg \min_{\beta} \hat{\sigma}(\mathbf{r}(\beta)) \quad (4.13)$$

where $\hat{\sigma}$ is a robust scale estimator.

Least Median of Squares (LMS) estimator:

The simplest robust scale estimator is the median of the absolute residuals:

$$\hat{\sigma}(\mathbf{r}) = \text{med}_i |r_i| \quad (4.14)$$

The corresponding regression estimator achieves highest breakdown point 50%, but a low efficiency. A tradeoff between breakdown and efficiency can be obtained by considering another quantile instead of the median.

Least Trimmed Squares (LTS) estimator:

Call $|r|_{(i)}$ the ordered absolute values of the residuals r_i , i.e. $|r|_{(1)} \leq \dots \leq |r|_{(n)}$. A smoother alternative is to consider a scale more similar to the standard deviation, namely the *trimmed squares scale*

$$\hat{\sigma}(\mathbf{r}) = \left(\frac{1}{h} \sum_{i=1}^h |r|_{(i)}^2 \right)^{1/2}, \quad (4.15)$$

with $h \in [n/2, n]$. Smaller values of h lead to higher breakdown point, but to lower efficiency.

M-estimator of scale

This estimator is defined as the solution σ of an equation of the form

$$\frac{1}{n} \sum_{i=1}^n \rho\left(\frac{r_i}{\hat{\sigma}}\right) = b \quad (4.16)$$

where ρ is a bounded ρ -function and b is a constant. It can be shown that the breakdown point of $\hat{\sigma}$ is $\min(b, 1 - b)$. Equation (4.16) is nonlinear, but it is easy to solve iteratively. Put

$$W(z) = \frac{\rho(z)}{z^2}. \quad (4.17)$$

Then (4.16) can be rewritten as

$$\hat{\sigma}^2 = \frac{1}{nb} \sum_{i=1}^n w_i r_i^2$$

with $w_i = W(r_i/\hat{\sigma})$, which displays $\hat{\sigma}$ as a weighted RMSE. Given some starting value $\hat{\sigma}_0$, an iterative procedure can be implemented as was done for regression M-estimators.

The choice $\rho(z) = z^2$ and $b = 1$ yields the RMSE. The choice $\rho(z) = I(|z| > 1)$ and $b = 0.5$ yields $\hat{\sigma} = \text{med}(|r|)$.

Regression estimators with $\hat{\sigma}$ given by (4.16) are called **S-estimators**. Although S-estimators achieve the maximum breakdown point, they have low efficiency. However, they can be used as initial estimator for the M-estimator (4.8). The resulting estimator is called **MM-estimator**; it inherits the breakdown point of the S-estimator, but has controllable efficiency.

14 Define affine equivariance.

4.4.1 Affine equivariance

It is desirable that location and covariance estimates respond in a mathematically convenient form to certain transformations of the data. One can define a transformation that is using a nonsingular $p \times p$ matrix \mathbf{A} and a vector \mathbf{b} of length p to transform the p -dimensional observations $\mathbf{x}_1, \dots, \mathbf{x}_n$ by $\mathbf{Ax}_j + \mathbf{b}$. This transformation performs any desired nonsingular linear transformation of the original data. Thus, if \mathbf{t} denotes a location estimator, it is requested that

$$\mathbf{t}(\mathbf{Ax}_1 + \mathbf{b}, \dots, \mathbf{Ax}_n + \mathbf{b}) = \mathbf{A} \cdot \mathbf{t}(\mathbf{x}_1, \dots, \mathbf{x}_n) + \mathbf{b}, \quad (4.18)$$

and for a covariance estimator \mathbf{C} we require

$$\mathbf{C}(\mathbf{Ax}_1 + \mathbf{b}, \dots, \mathbf{Ax}_n + \mathbf{b}) = \mathbf{A} \cdot \mathbf{C}(\mathbf{x}_1, \dots, \mathbf{x}_n) \cdot \mathbf{A}^\top. \quad (4.19)$$

Location and covariance estimators that fulfil (4.18) and (4.19) are called *affine equivariant* estimators. These estimators transform properly under changes of the origin, the scale, or under rotations.

The property of affine equivariance is only valid for nonsingular transformation matrices. Note that the coordinate-wise median as a robust location estimator is not affine equivariant. Also a robust covariance estimate based on pairwise robust covariances would not be affine equivariant.

15 What is the Minimum Covariance Determinant estimator? How does it work (in concept)? What properties does it have (related to tuning parameter h)?

An estimator of multivariate location and covariance which is affine equivariant and has high breakdown point is the Minimum Covariance

Determinant (MCD) estimator. The idea behind this estimator is in fact related to the LTS estimator. Here, one is searching for those h data points for which the determinant of the empirical covariance matrix is minimal. The location estimator \mathbf{t} is the mean of these h observations, and the covariance estimator \mathbf{C} is given by the covariance matrix with the smallest determinant, but multiplied by a constant to obtain consistency for normal distribution. The parameter h determines the robustness but also the efficiency of the resulting estimator. The highest possible breakdown point can be achieved if $h \approx n/2$ is taken, but this choice leads to a low efficiency. On the other hand, for higher values of h the efficiency increases but the breakdown point decreases. Therefore, a compromise between efficiency and robustness is considered in practice. The data coming from the outlier distribution are inflating the tolerance ellipse based on the classical estimators while that based on the MCD is much more compact and reflects the structure of the majority of data.

Multivariate S-estimators

Similar as in the regression context, see Equation (4.16), it is possible to define S-estimators in the context of robust location and covariance estimation. The idea is to make the Mahalanobis distances small. The squared Mahalanobis distances are defined as

$$\text{MD}^2(\mathbf{x}_i, \mathbf{t}, \mathbf{C}) = (\mathbf{x}_i - \mathbf{t})^\top \mathbf{C}^{-1} (\mathbf{x}_i - \mathbf{t}) \quad \text{for } i = 1, \dots, n$$

for a location estimator \mathbf{t} and a covariance estimator \mathbf{C} .

Small Mahalanobis distances can be achieved by using an M-estimator of scale $\hat{\sigma}$, and minimizing

$$\hat{\sigma}(\text{MD}^2(\mathbf{x}_1, \mathbf{t}, \mathbf{C}), \dots, \text{MD}^2(\mathbf{x}_n, \mathbf{t}, \mathbf{C}))$$

under the restriction that the determinant $|\mathbf{C}| = 1$. This restriction avoids a degenerated solution for \mathbf{C} .

Multivariate MM-estimators

Like in robust regression, a drawback of S-estimators is that their asymptotic efficiency might be rather low. MM estimators for multivariate location and covariance combine both high breakdown point and high efficiency. The resulting estimators are affine equivariant and have bounded influence function. As for the previous estimators, the solution for the estimators can be found by an iterative algorithm.

- 16 What are problems in classic regression diagnostics (hat matrix)? What are robust regression diagnostics?**

LS regression is sensitive to outliers. These could be outliers in the response, so-called “vertical” outliers, or outliers in the explanatory variables, so-called “leverage points”. The latter are particularly influential if they have large residual (bad leverage points).

In LS regression, the hat matrix \mathbf{H} is used to identify leverage points. The hat matrix is defined as $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$, and the LS fit is obtained as $\hat{\mathbf{y}}_{LS} = \mathbf{H}\mathbf{y}$.

If there is an exact fit, one can assume that this was caused by a leverage point, because an extremely strong influence was caused to the own estimation. In general, one needs to be careful if there are large values of h_{ii} . As a rule of thumb one could define the threshold $h_{ii} > 2 \cdot \frac{p+1}{n}$ for the identification of leverage points.

What is the reason for the masking effect? In fact, one can show that the following relationship holds in case of regression with intercept:

$$h_{ii} = \frac{1}{n-1} \text{MD}_i^2 + \frac{1}{n} \quad (4.20)$$

Here,

$$\text{MD}_i^2 = (\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{S}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}})$$

Regression diagnostics

The different types of outliers we would like to distinguish are illustrated in the case of simple linear regression schematically in Figure 4.9. These are:

- Regular observations: \mathbf{x}_i is in the usual data range, and y_i fits to the model.
- Vertical outliers: \mathbf{x}_i is in the usual data range, but y_i does not fit to the model.
- Good leverage points: \mathbf{x}_i is an outlier, thus unusual in the x-space, but y_i fits to the model.

- Bad leverage points. x_i is an outlier, thus unusual in the x-space, and y_i does not fit to the model.

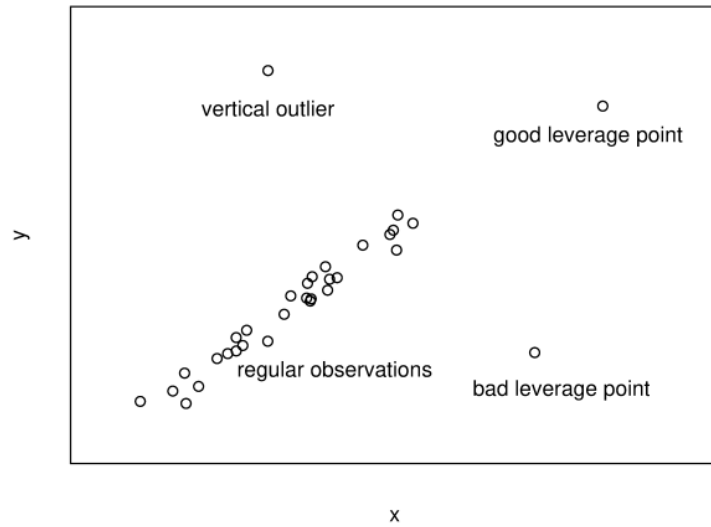


Figure 4.9: Different types of outliers, illustrated in simple linear regression.

In general, good leverage points have the advantage that they are along the regression line (hyperplane) and thus they even allow for a more accurate estimation of the regression parameters. Bad leverage points can have a strong impact on the (LS) estimation, and they can even lead to a leverage of the regression line (hyperplane).

The regression diagnostic plot allows to distinguish these 4 types of observations. Outlyingness in the x-space can be recognized by robust Mahalanobis distances. Large (absolute) residuals can be recognized as scaled deviations in the y-space from the robust regression fit. Thus, the robust residuals need to be scaled by a

robust estimate of the residual scale. Then one can argue according to the normal theory that scaled residuals outside $[-2, 2]$ or $[-2.5, 2.5]$ are extreme and thus very unusual.

Figure 4.10 shows the resulting regression diagnostic plot, where the robust Mahalanobis distances are plotted against the robust scaled residuals. The horizontal lines are at the thresholds ± 2.5 for the scaled residuals, and the vertical line is at the threshold based on the chi-square quantile.

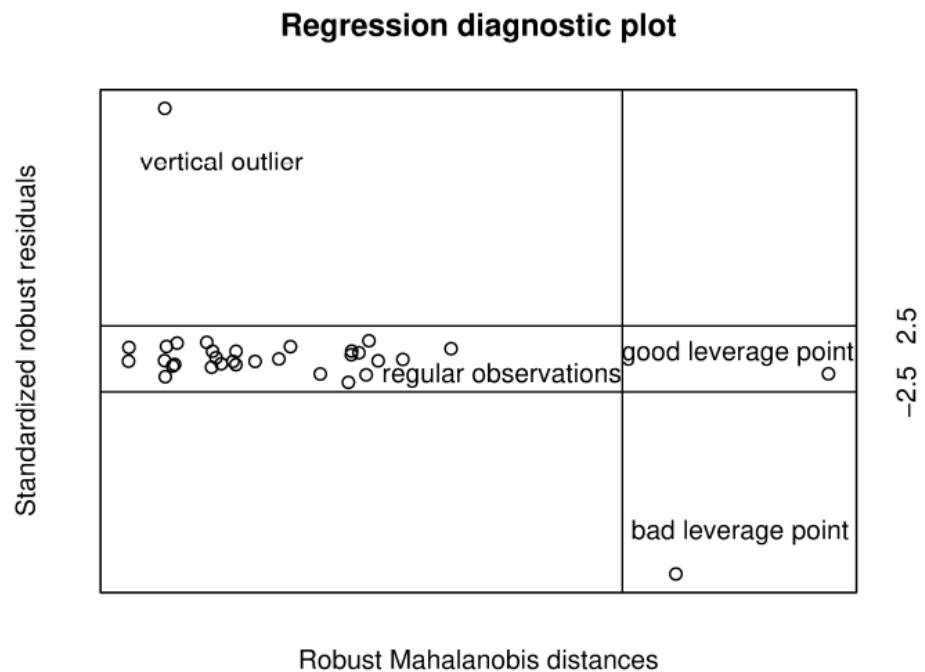


Figure 4.10: Identification of the four different categories of observations in the regression diagnostic plot.

17 How does robust multivariate regression work? (estimate covariance matrix with M-estimator of scale)

4.6 Robust multivariate regression

Consider again the multivariate regression problem with the observations $\mathbf{y}_{i.} = (y_{i1}, \dots, y_{im})^\top$, $\mathbf{x}_{i.} = (1, x_{i1}, \dots, x_{ip})^\top$, and the error terms $\mathbf{e}_{i.} = (e_{i1}, \dots, e_{im})^\top$, and the model

$$\mathbf{y}_{i.} = \mathbf{B}^\top \mathbf{x}_{i.} + \mathbf{e}_{i.} ,$$

for $i = 1, \dots, n$, see also Equation (3.13). We are interested in robustly estimating the $(p+1) \times m$ matrix \mathbf{B} of regression coefficients, as well as the covariance matrix Σ of the error terms. For this purpose, one can use multivariate S-estimators, see Section 4.4.3, and Van Aelst and Willems (2005). Here, the Mahalanobis distances are based on the residuals $\mathbf{r}_{i.} = \mathbf{y}_{i.} - \mathbf{B}^\top \mathbf{x}_{i.}$ and on an estimator \mathbf{C} of Σ . Thus, one has to minimize

$$\hat{\sigma}(\mathbf{r}_{1.}^\top \mathbf{C}^{-1} \mathbf{r}_{1.}, \dots, \mathbf{r}_{n.}^\top \mathbf{C}^{-1} \mathbf{r}_{n.})$$

using an M-estimator of scale $\hat{\sigma}$, under the constraint $|\mathbf{C}| = 1$.

18 Principal Component Analysis - how to select the vectors for the transformation, Lagrange problem definition.

Aim is to describe complex relationships in given data in a simpler form. The data are represented by linear combinations of specific components in a way to preserve as much information as possible. Thus, the dimensionality is reduced to the number of those components.

5.2 Definition of principal components

Let $\mathbf{x} = (x_1, \dots, x_p)^\top$ be a p -dimensional random vector with expectation $E(\mathbf{x}) = \boldsymbol{\mu}$ and covariance matrix

$$\Sigma = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] . \quad (5.1)$$

Further, $\mathbf{\Gamma} = (\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_p)$ is a $(p \times p)$ matrix with fixed values (non-random), with the constraint that its column vectors $\boldsymbol{\gamma}_i$ are unitary vectors, i.e. $\boldsymbol{\gamma}_i^\top \boldsymbol{\gamma}_i = 1$ for $i = 1, \dots, p$. Moreover, different columns of $\mathbf{\Gamma}$ are orthogonal, i.e. $\boldsymbol{\gamma}_i^\top \boldsymbol{\gamma}_j = 0$ for $i \neq j$. This implies that $\mathbf{\Gamma}^\top = \mathbf{\Gamma}^{-1}$.

Consider the linear transformation

$$\mathbf{z} = \mathbf{\Gamma}^\top (\mathbf{x} - \boldsymbol{\mu}) \quad (5.2)$$

or, expressed in components,

$$z_i = \boldsymbol{\gamma}_i^\top (\mathbf{x} - \boldsymbol{\mu}) \quad \text{for} \quad i = 1, \dots, p . \quad (5.3)$$

The result of the above transformation is a new random variable \mathbf{z} of dimension p . The variance of z_i ($i = 1, \dots, p$) is

$$\text{Var}(z_i) = E\left[\gamma_i^\top (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \gamma_i\right] = \gamma_i^\top \boldsymbol{\Sigma} \gamma_i . \quad (5.4)$$

So far it is not clear which matrix $\boldsymbol{\Gamma}$ we should choose – in spite of the constraints, there are still infinitely many possibilities. Since we are interested in preserving information content, we should look at the variance of the transformed variables. In other words, we would like to obtain such a transformation which maximizes the variances of the components of \mathbf{z} .

Considering also the constraints on $\boldsymbol{\Gamma}$, we can mathematically formulate a maximization problem in terms of Lagrange optimization.

First principal component:

For $i = 1$ we would like to obtain a component z_1 such that $\text{Var}(z_1)$ is maximized, and $\gamma_1^\top \gamma_1 = 1$. The corresponding Lagrangian problem can be stated as:

$$\phi_1 = \gamma_1^\top \boldsymbol{\Sigma} \gamma_1 - a_1(\gamma_1^\top \gamma_1 - 1) \quad (5.5)$$

The partial derivatives with respect to the unknowns γ_1 are set equal to zero, and we obtain

$$\frac{\partial \phi_1}{\partial \gamma_1} = 2\boldsymbol{\Sigma} \gamma_1 - 2a_1 \gamma_1 = \mathbf{0} \quad (5.6)$$

or

$$\boldsymbol{\Sigma} \gamma_1 = a_1 \gamma_1 . \quad (5.7)$$

This is an eigenvector/eigenvalue problem, and the solution is that the unknown coefficient vector γ_1 is an eigenvector of the covariance matrix $\boldsymbol{\Sigma}$ to the eigenvalue a_1 .

Our problem, however, is that $\boldsymbol{\Sigma}$ has p eigenvectors in total. Which one should we take?

With Equation (5.7) we can see that

$$\text{Var}(z_1) = \gamma_1^\top (\boldsymbol{\Sigma} \gamma_1) = \gamma_1^\top (a_1 \gamma_1) = a_1 \gamma_1^\top \gamma_1 = a_1,$$

and since we want to maximize variance, we take that eigenvector γ_1 corresponding to the largest eigenvalue a_1 . Component z_1 now is denoted as *first principal component (PC)*, and γ_1 is the direction of this component.

Second principal component:

In the next step, for $i = 2$, we again want to maximize variance, more clearly, $\text{Var}(z_2)$ should be maximized under the constraint $\gamma_2^\top \gamma_2 = 1$. In addition we also want that z_1 and z_2 are uncorrelated (in order to uncover new information in z_2). The latter condition means:

$$\begin{aligned} \text{Cov}(z_1, z_2) &= \text{Cov}(\gamma_1^\top (\mathbf{x} - \boldsymbol{\mu}), \gamma_2^\top (\mathbf{x} - \boldsymbol{\mu})) = \text{Cov}(\gamma_1^\top \mathbf{x}, \gamma_2^\top \mathbf{x}) = \gamma_1^\top \boldsymbol{\Sigma} \gamma_2 = \\ &= \gamma_2^\top \boldsymbol{\Sigma} \gamma_1 = \gamma_2^\top a_1 \gamma_1 = a_1 \gamma_2^\top \gamma_1 = 0 \end{aligned}$$

Since $a_1 \neq 0$, the condition of uncorrelated components is equivalent to orthogonality of γ_1 and γ_2 . That's a remarkable fact which you don't easily get for other methods.

Now we can again formulate the Lagrangian problem,

$$\phi_2 = \gamma_2^\top \Sigma \gamma_2 - a_2(\gamma_2^\top \gamma_2 - 1) - b\gamma_2^\top \gamma_1 \quad (5.8)$$

with the Lagrange coefficients a_2 and b . The partial derivatives with respect to the unknowns γ_2 are set equal to zero, which yields:

$$\frac{\partial \phi_2}{\partial \gamma_2} = 2\Sigma \gamma_2 - 2a_2\gamma_2 - b\gamma_1 = \mathbf{0} \quad (5.9)$$

Multiplication from the left-hand side with γ_1^\top result in

$$2\gamma_1^\top \Sigma \gamma_2 - 2a_2\gamma_1^\top \gamma_2 - b\gamma_1^\top \gamma_1 = 0 - 0 - b \cdot 1 = 0 ,$$

and thus $b = 0$. Therefore we can reduce (5.9) to

$$\Sigma \gamma_2 = a_2 \gamma_2 ,$$

and thus γ_2 is eigenvector to Σ to the next largest eigenvalue a_2 , and z_2 is called *second PC*.

19 What is the expectation of the Principal Components?

5.4 PCs based on data

Consider an $n \times p$ data matrix \mathbf{X} with n observations and p variables. For the PCA transformation we need to estimate the expectation vector and the covariance matrix. This can be done by the empirical estimates, the sample mean and the sample covariance matrix,

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad (5.20)$$

and

$$\hat{\boldsymbol{\Sigma}} = \mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top, \quad (5.21)$$

where \mathbf{x}_i is the i -th row of \mathbf{X} .

In analogy to Equation (5.11), the sample principal components are computed by

$$\mathbf{Z} = (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^\top)\hat{\boldsymbol{\Gamma}} \quad (5.22)$$

or

$$\mathbf{z}_j = (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^\top)\hat{\boldsymbol{\gamma}}_j \quad \text{for} \quad j = 1, \dots, p. \quad (5.23)$$

Here, $\mathbf{1}$ is a column vector with n entries of 1, and $\hat{\boldsymbol{\Gamma}} = (\hat{\boldsymbol{\gamma}}_1, \dots, \hat{\boldsymbol{\gamma}}_p)$ is the matrix of eigenvectors of \mathbf{S} . Similar to the population definition of PCA, we obtain a decomposition

$$\hat{\boldsymbol{\Gamma}}^\top \mathbf{S} \hat{\boldsymbol{\Gamma}} = \hat{\mathbf{A}}, \quad (5.24)$$

where $\hat{\mathbf{A}} = \text{Diag}(\hat{a}_1, \dots, \hat{a}_p)$ is a diagonal matrix with the eigenvalues to the corresponding eigenvectors of \mathbf{S} , arranged in descending order.

The matrix \mathbf{Z} of PCs has the same dimension as the data matrix \mathbf{X} . In fact, it just represents the data information in an orthogonally rotated coordinate system.

In this representation, the data are centered, and thus the original data center could not be reconstructed. The elements z_{ij} of the matrix \mathbf{Z} are called (PC) *scores*.

In analogy to Equation (5.17) we can compute the correlations between the variables and the PCs as

$$\hat{\lambda}_{ij} = \frac{\hat{\gamma}_{ij}\hat{a}_j^{\frac{1}{2}}}{s_{ii}^{\frac{1}{2}}} \quad \text{for} \quad i, j = 1, \dots, p, \quad (5.25)$$

where s_{ii} is the i -th diagonal element of \mathbf{S} . In matrix notation, this correlation is

$$\hat{\boldsymbol{\Lambda}} = \left(\text{Diag}(\mathbf{S}) \right)^{-\frac{1}{2}} \hat{\boldsymbol{\Gamma}} \hat{\mathbf{A}}^{\frac{1}{2}}. \quad (5.26)$$

The contribution of the variables to the PCs are the loadings.

The empirical variances of the PCs are equal to the eigenvalues.

20 Why is PCA sensitive to scale? What happens if we center-scale the data?

Remark: We can see that the PC transformation is not scale invariant, which means that the resulting PCs depend on the scale (units) of the variables. If we would first scale (and probably also center) the variables to mean 0 and variance 1,

$$y_i = \frac{x_i - E(x_i)}{\sqrt{Var(x_i)}} \quad \text{for } i = 1, \dots, p, \quad (5.12)$$

we would in general obtain different PCs, since then we would perform the eigenvector/eigenvalue decomposition on the correlation matrix, and not on the covariance matrix. Therefore, if scale invariance is desired (and in most applications it is), the variables need to be standardized first.

Since our random variable \mathbf{x} has been centered, we obtain that the expectation of the PCs is zero:

$$E(\mathbf{z}) = \mathbf{\Gamma}^\top [E(\mathbf{x} - \boldsymbol{\mu})] = \mathbf{0} \quad (5.13)$$

The covariance matrix of the PCs is

$$\text{Cov}(\mathbf{z}) = \mathbf{\Gamma}^\top \text{Cov}(\mathbf{x} - \boldsymbol{\mu}) \mathbf{\Gamma} = \mathbf{\Gamma}^\top \mathbf{\Sigma} \mathbf{\Gamma} = \mathbf{A} . \quad (5.14)$$

21 What are some rules for the number of principal components to select? (for hypothesis tests: only concept, not formulas)

The answer might depend on the purpose, what the user wants to do with the PCs. If it is for visual inspection of the data, it might be sufficient to look at those PCs which are covering the most important data information. One could also argue, that part of the information just consists of noise, and this should be contained in the last few PCs, which are not interesting for the inspection.

PCA tries to explain as much of the total variance as possible with fewer dimensions. What is the total variance? It is:

$$\sum_{j=1}^p s_{jj} = \sum_{j=1}^p \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$$

This is the same as,

$$\frac{1}{n-1} \text{trace}((\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^\top)^\top (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^\top)) = \text{trace}(\mathbf{S}) = \text{trace}(\hat{\mathbf{\Gamma}} \hat{\mathbf{A}} \hat{\mathbf{\Gamma}}^\top) = \text{trace}(\hat{\mathbf{A}}) ,$$

where “trace” is the sum of the diagonal elements of the matrix, here the sum of the eigenvalues.

There exist several other tests, such as a test of the hypothesis that the explained variance of the first k PCs exceeds a certain threshold, such as 80 % or 90 %.

A further frequently applied criterion is to exclude those PCs which have a variance (eigenvalue) lower than the average. If the data are standardized, the sum of the eigenvalues is p , and thus the average is 1.

A further possibility to select the number of relevant PCs is the *scree graph*, which shows the proportion of explained variance of each PC versus the number of the PC. We exclude those (smallest) PCs where the proportion follows approximately a linear trend.

Example 5.5.3 The scree graph for the exam data is shown in Figure 5.4. One can see that the contributions of the last 3 PCs follow a linear trend, and thus these PCs can be excluded.

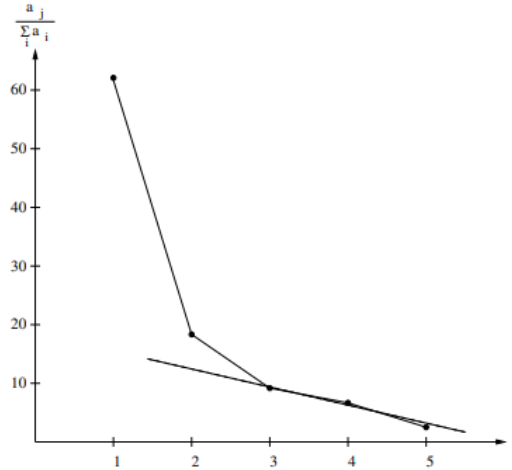


Figure 5.4: Scree graph for the exam data example.

- 22 **What is singular value decomposition, how is it defined, and how is it related to PCA? What are the scores in terms of SVD? When would we prefer SVD to spectral decomposition of the covariance (correlation) matrix?**

5.6 Singular value decomposition

Singular value decomposition (SVD) can be viewed as an alternative algorithm to determine the PCs. It is not based on a decomposition of the covariance matrix, but it directly uses the data matrix for the decomposition. This is a particular advantage if $n < p$ (“flat” data matrices), which would always lead to a singular covariance matrix, where the last $n - p$ PCs would have variance 0.

Let us assume in the following that the columns of the real-valued data matrix \mathbf{X} are centered to mean 0. Then there exists an orthogonal $n \times n$ matrix \mathbf{U} and an orthogonal $p \times p$ matrix \mathbf{V} , such that we obtain the decomposition

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top, \quad (5.31)$$

where \mathbf{D} is an $n \times p$ matrix with “diagonal” elements $d_{ii} \geq 0$ for $i = 1, \dots, \min(n, p)$, and the remaining elements are 0. The positive values d_{ii} are called *singular values* of \mathbf{X} . The number of these positive values corresponds to the rank of \mathbf{X} .

If the rank of \mathbf{X} is $k \leq \min(n, p)$, then \mathbf{X} can be represented as

$$\mathbf{X} = \sum_{i=1}^k d_{ii} \mathbf{u}_i \mathbf{v}_i^\top, \quad (5.32)$$

where \mathbf{u}_i and \mathbf{v}_i are the i -th columns of \mathbf{U} and \mathbf{V} , respectively. Due to the orthogonality of \mathbf{U} and \mathbf{V} we obtain

$$\mathbf{X} \mathbf{X}^\top \mathbf{u}_i = d_{ii}^2 \mathbf{u}_i \quad (5.33)$$

and

$$\mathbf{X}^\top \mathbf{X} \mathbf{v}_i = d_{ii}^2 \mathbf{v}_i. \quad (5.34)$$

This means that \mathbf{u}_i is the i -th eigenvector of $\mathbf{X}\mathbf{X}^\top$ to the eigenvalue d_{ii}^2 , and \mathbf{v}_i is the i -th eigenvector of $\mathbf{X}^\top\mathbf{X}$ to the same eigenvalue d_{ii}^2 . The eigenvalues for $i = 1, \dots, k$ are strictly positive, and the remaining ones are zero.

As a summary we conclude that \mathbf{U} has n orthogonal eigenvectors of $\mathbf{X}\mathbf{X}^\top$ in its columns, and \mathbf{V} has p orthogonal eigenvectors of $\mathbf{X}^\top\mathbf{X}$ in its columns.

It is now straightforward to see the connection to a covariance-based estimation of the PCs, as it was done in Section 5.4: Consider again mean-centered data \mathbf{X} .

Then the sample PCs were defined as $\mathbf{Z} = \mathbf{X}\mathbf{\Gamma}$, and thus $\mathbf{X} = \mathbf{Z}\mathbf{\Gamma}^\top$. Further, in this case the sample covariance matrix is

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}^\top \mathbf{X} = \hat{\mathbf{\Gamma}} \hat{\mathbf{A}} \hat{\mathbf{\Gamma}}^\top.$$

The matrix $\hat{\mathbf{\Gamma}}$ is the matrix with the normed eigenvectors of \mathbf{S} . In SVD, \mathbf{V} is the matrix with the normed eigenvectors of $\mathbf{X}^\top\mathbf{X}$, and therefore we can conclude that $\hat{\mathbf{\Gamma}} \equiv \mathbf{V}$. From Equation (5.34) we see that $d_{ii}^2 = (n-1)\hat{a}_i$. We can thus write

$$\mathbf{X} = \mathbf{Z}\mathbf{V}^\top = \mathbf{U}\mathbf{D}\mathbf{V}^\top, \quad (5.35)$$

and therefore we have $\mathbf{Z} = \mathbf{U}\mathbf{D}$.

23 How can we define the PCA problem in terms of reconstruction error (Frobenius norm)?

Alternative definitions for PCs

With these considerations we can formulate the PCA problem based on a different objective function.

Let us first define the *Frobenius norm* of a matrix: Denote \mathbf{x}_i as the rows of \mathbf{X} , for $i = 1, \dots, n$. The Frobenius norm of \mathbf{X} is defined as:

$$\|\mathbf{X}\|_F = \sqrt{\sum_{i=1}^n \|\mathbf{x}_i\|_2^2} = \sqrt{\sum_{i=1}^n \sum_{j=1}^p x_{ij}^2}$$

From above we have

$$\mathbf{X}\mathbf{V} = \mathbf{U}\mathbf{D} = \mathbf{Z}.$$

Define $\mathbf{V}_m = (\mathbf{v}_1, \dots, \mathbf{v}_m)$, and m is smaller than the rank of \mathbf{X} . Then the columns of $\mathbf{X}\mathbf{V}_m$ are the first m PCs. This is equivalent to a projection of \mathbf{X} onto an m -dimensional subspace formed by \mathbf{V}_m . One can show that

$$\mathbf{V}_m = \arg \max_{\mathbf{B}} \|\mathbf{X}\mathbf{B}\|_F^2$$

for any $p \times m$ matrix \mathbf{B} with $\text{rank}(\mathbf{B}) \leq m$ and $\mathbf{B}^\top \mathbf{B} = \mathbf{I}$.

An equivalent formulation is the following: We have

$$\mathbf{X} = \mathbf{X}\mathbf{V}\mathbf{V}^\top = \mathbf{X}\mathbf{V}_m\mathbf{V}_m^\top + \mathbf{E}.$$

Then

$$\mathbf{V}_m = \arg \min_{\mathbf{B}} \|\mathbf{X} - \mathbf{X}\mathbf{B}\mathbf{B}^\top\|_F^2$$

with the same definition of \mathbf{B} as above. Note that $\hat{\mathbf{X}} := \mathbf{X}\mathbf{B}\mathbf{B}^\top$ has the same dimension as \mathbf{X} . Therefore we can view $\hat{\mathbf{X}}$ as a rank m approximation of \mathbf{X} , which is optimal in the above sense. For finding the PC directions we are in fact minimizing residual sum-of-squares, with the residual matrix $\mathbf{X} - \hat{\mathbf{X}}$.

24 What are Biplots? What is the rank-2 approximation? Define the G/H matrix. What are the properties of the biplot? (inner row product of G and H approximates elements of the X-matrix, etc.)

Biplots have been introduced in Gabriel (1971) to display both variable and object information jointly in one plot – usually in 2 dimensions. The “bi” does not refer to the “2 dimensions” but to the joint presentation of variables and observations. Here we will introduce biplots to represent loadings and scores from a PCA.

We like 2-dimensional plots because they are easy to handle. A projection of the data into the plane should, however, make sure that the data have approximately rank 2. Then we are sure that indeed the projection represents the essential variability.

Let \mathbf{X} be a mean-centered $n \times p$ data matrix of rank k (not much larger than 2). The biplot shows a representation of \mathbf{X} by means of two groups of vectors with dimension n and p , respectively, which form a rank-2 approximation of \mathbf{X} . Denote this rank-2 approximation by $\mathbf{X}_{(2)}$.

We know that a least-squares based rank-2 approximation is given by the first 2 PCs, and we are aware of SVD to compute these first 2 PCs. In the following we will only need the first 2 columns of \mathbf{U} and \mathbf{V} , but we will still use the notation \mathbf{U} and \mathbf{V} in order to avoid new symbols. Thus:

$$\mathbf{X} \approx \mathbf{X}_{(2)} = \mathbf{U} \mathbf{D} \mathbf{V}^\top = (\mathbf{u}_1 \ \mathbf{u}_2) \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix} \begin{pmatrix} \mathbf{v}_1^\top \\ \mathbf{v}_2^\top \end{pmatrix}. \quad (5.36)$$

Since we want to represent the data by n score vectors and p loadings vectors (all 2-dimensional), we need a decomposition of $\mathbf{X}_{(2)}$ into an $n \times 2$ and an $p \times 2$ matrix, say

$$\mathbf{X}_{(2)} = \mathbf{G} \mathbf{H}^\top \quad (5.37)$$

with

$$\mathbf{G} = (\mathbf{u}_1 \ \mathbf{u}_2) \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix}^{1-c} \quad (5.38)$$

and

$$\mathbf{H} = (\mathbf{v}_1 \ \mathbf{v}_2) \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix}^c \quad (5.39)$$

for $0 \leq c \leq 1$. Depending on the choice of c , the first 2 singular values are distributed among the matrices \mathbf{G} and \mathbf{H} . The biplot consists of the rows of \mathbf{G} and \mathbf{H} , i.e. of $n + p$ 2-dimensional vectors.

For the choice $c = 0.5$ we use the same scaling for the observation and variable vectors, which seems to be natural. However, we obtain nicer properties with the

choice $c = 1$ (and with re-scaling):

$$\mathbf{G} = \begin{pmatrix} \mathbf{g}_{1.}^\top \\ \vdots \\ \mathbf{g}_{n.}^\top \end{pmatrix} = \sqrt{n-1} \begin{pmatrix} \mathbf{u}_1 & \mathbf{u}_2 \end{pmatrix} \quad (5.40)$$

$$\mathbf{H} = \begin{pmatrix} \mathbf{h}_{1.}^\top \\ \vdots \\ \mathbf{h}_{p.}^\top \end{pmatrix} = \frac{1}{\sqrt{n-1}} \begin{pmatrix} \mathbf{v}_1 & \mathbf{v}_2 \end{pmatrix} \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix} \quad (5.41)$$

These properties are:

- The inner product between the rows of \mathbf{G} and the rows of \mathbf{H} approximate the values x_{ij} of the data matrix \mathbf{X} :

$$\mathbf{g}_{i.}^\top \mathbf{h}_{j.} = \sqrt{n-1} \mathbf{u}_{i.}^\top \frac{1}{\sqrt{n-1}} (\mathbf{v}_{j.}^\top \mathbf{D})^\top = \mathbf{u}_{i.}^\top \mathbf{D} \mathbf{v}_{j.} \approx x_{ij} . \quad (5.42)$$

- The inner product between the rows of \mathbf{H} approximates the covariance:

$$\begin{aligned} \mathbf{H} \mathbf{H}^\top &= \left(\frac{1}{\sqrt{n-1}} \mathbf{V} \mathbf{D} \right) \left(\frac{1}{\sqrt{n-1}} \mathbf{D} \mathbf{V}^\top \right) = \frac{1}{n-1} \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top \\ &= \frac{1}{n-1} (\mathbf{V} \mathbf{D} \mathbf{U}^\top) (\mathbf{U} \mathbf{D} \mathbf{V}^\top) = \frac{1}{n-1} \mathbf{X}_{(2)}^\top \mathbf{X}_{(2)} \\ &\approx \frac{1}{n-1} \mathbf{X}^\top \mathbf{X} = \mathbf{S} . \end{aligned} \quad (5.43)$$

It follows that the squared Euclidean norm of the rows of \mathbf{H} , i.e. $\|\mathbf{h}_{j.}\|^2 = \mathbf{h}_{j.}^\top \mathbf{h}_{j.}$, approximates the variance. Moreover, the cosine between $\mathbf{h}_{i.}$ and $\mathbf{h}_{j.}$ ($i, j = 1, \dots, p$) approximates the correlation between the variables,

$$\cos(\mathbf{h}_{i.}, \mathbf{h}_{j.}) = \frac{\mathbf{h}_{i.}^\top \mathbf{h}_{j.}}{\|\mathbf{h}_{i.}\| \|\mathbf{h}_{j.}\|} \approx r_{ij} . \quad (5.44)$$

- The Euclidean distance between the rows of \mathbf{G} approximates the Mahalanobis distance between the observations,

$$\begin{aligned} \|\mathbf{g}_{i.} - \mathbf{g}_{j.}\|^2 &= (\mathbf{g}_{i.} - \mathbf{g}_{j.})^\top (\mathbf{g}_{i.} - \mathbf{g}_{j.}) = (n-1) (\mathbf{u}_{i.} - \mathbf{u}_{j.})^\top (\mathbf{u}_{i.} - \mathbf{u}_{j.}) \\ &\approx (\mathbf{x}_{i.} - \mathbf{x}_{j.})^\top \mathbf{S}^{-1} (\mathbf{x}_{i.} - \mathbf{x}_{j.}) , \end{aligned} \quad (5.45)$$

because

$$\mathbf{x}_{i.}^\top \mathbf{S}^{-1} \mathbf{x}_{j.} \approx (\mathbf{u}_{i.}^\top \mathbf{D} \mathbf{V}^\top) (n-1) (\mathbf{V} \mathbf{D}^{-2} \mathbf{V}^\top) (\mathbf{V} \mathbf{D} \mathbf{u}_{j.}) = (n-1) \mathbf{u}_{i.}^\top \mathbf{u}_{j.} \quad (5.46)$$

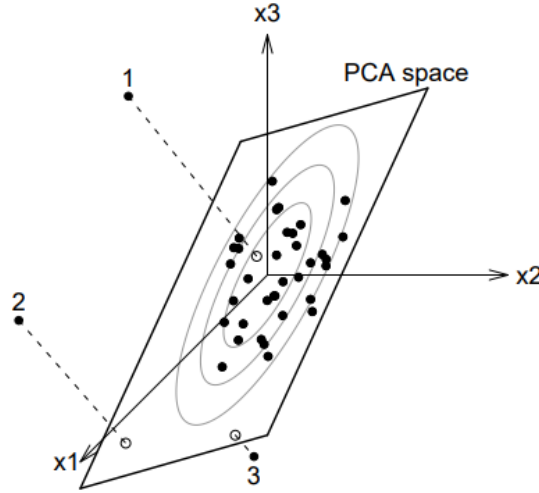
for $i, j = 1, \dots, n$.

25 Which diagnostics do we have for PCA (formal definition of orthogonal, score distance)?

Two distance measures have been introduced: the score distance (SD) and the orthogonal distance (OD).

Figure 5.6 tries to explain the ideas behind these distances. We have observations in the 3-dimensional space, and the PCA space is built up with 2 PCs. The SD is a distance measure in the PCA space, and it is equal to the Mahalanobis distances of the observations projected into this space. The ellipses refer to these Mahalanobis distances. The OD is the distance orthogonal to the PCA space.

Here we can see three particular observations: Observation 1 has big OD but small SD; observation 2 has big OD and big SD; observation 3 has small OD but big SD. Similar to the diagnostics in regression one could classify these different types of observations as vertical outliers (1), and good (3) and bad (2) leverage points. In particular, bad leverage points can have a strong influence on the classical estimation of the PCs. Here, “classical” refers to the SVD estimation, or equivalently, the estimation based on the empirical covariance matrix. A robust PCA could easily be obtained through a robust estimation of the covariance matrix. In this case, the diagnostics would also make much more sense.



Formally, these distances are defined as follows. Let $\hat{\mathbf{\Gamma}}_k = (\hat{\gamma}_1, \dots, \hat{\gamma}_k)$ be the matrix of the first k estimated PC loadings, $\mathbf{z}_i = (z_{i1}, \dots, z_{ik})^\top$ the i -th score vector, and $\hat{a}_1, \dots, \hat{a}_k$ the corresponding variances of the PCs. The number k of PCs can be selected according to a criterion defined in Section 5.5, e.g. such that 80% of the variability is explained.

The SD for the i -th observation ($i = 1, \dots, n$) is defined as

$$SD_i = \left(\sum_{j=1}^k \frac{z_{ij}^2}{\hat{a}_j} \right)^{1/2}.$$

This is equal to the Mahalanobis distance of the score vector to the PCA center (zero) with respect to the covariance matrix (\mathbf{A}).

The OD for the i -th observation \mathbf{x}_i ($i = 1, \dots, n$) is defined as

$$OD_i = \|\mathbf{x}_i - \hat{\mathbf{\Gamma}}_k \mathbf{z}_i\|_2,$$

which is the Euclidean distance of the observation to its projection into the space of the first k PCs.

Similar to multivariate outlier detection, one can define cutoff values for both distance measures, which would refer to unusual values of SD and/or OD. Since SD is a Mahalanobis distance, a suitable cutoff value is $\sqrt{\chi_{k;0.975}^2}$. For the cutoff value for the OD it has been argued that $OD^{2/3}$ is closer to normality, and thus a suitable cutoff value is

$$\left(\text{median}_i(OD_i^{2/3}) + \text{MAD}_i(OD_i^{2/3}) z_{0.975} \right)^{3/2},$$

where $z_{0.975}$ is the 0.975 quantile of the $N(0, 1)$. MAD stands for the Median Absolute Deviation, which is defined for univariate values y_1, \dots, y_n as

$$\text{MAD} = 1.483 \cdot \text{median}_i(|y_i - \text{median}_j(y_j)|).$$

26 What is the factor analysis model (formal definition, assumptions)? What is the difference to PCA?

In factor analysis we assume that what we observe is basically the result of underlying quantities which are not directly observable. These quantities are called latent variables, and they cannot be measured. The “factors” in factor analysis aim at isolating such latent variables, explaining the relationships in the data.

Note that factor analysis is basically similar to PCA, since we also aim for a dimension reduction. However, we would like to have interpretable factors (the components in PCA are not necessarily interpretable), and we also have a statistical model (PCs were only defined by a linear transformation).

27 Explain the decomposition of the correlation matrix in factor analysis.

6.2 Factor analysis model

6.2.1 Definition

Let $\mathbf{x} = (x_1, x_2, \dots, x_p)^\top$ be a p -dimensional vector of random variables, x_1, x_2, \dots, x_p , which will describe our characteristics or variables later on.

Since in PCA we had the problem of the dependency on the scale, we will right away center and scale the random variables to mean zero and variance one. Thus, $\mathbf{y} = (y_1, y_2, \dots, y_p)^\top$ is the new random variable, obtained by

$$y_i = \frac{x_i - E(x_i)}{\sqrt{\text{Var}(x_i)}}, \quad i = 1, \dots, p.$$

In the model we already assume that the information contained in \mathbf{y} can be re-expressed by a smaller number $k < p$ of unknown random variables (factors) $\mathbf{f} = (f_1, \dots, f_k)^\top$, up to an error term \mathbf{e} . Thus, dimension reduction is already intrinsically stated in the model, called k -factor model:

$$\mathbf{y} = \mathbf{\Lambda} \mathbf{f} + \mathbf{e}. \quad (6.1)$$

Here, $\mathbf{\Lambda} = [(\lambda_{ij})]$ is a $(p \times k)$ matrix with fixed values. It is called *loadings matrix*, and it describes the relationships between factors and variables (as in PCA). The error term $\mathbf{e} = (e_1, \dots, e_p)^\top$ is often called *unique factor* (or uniqueness).

In this model we have the following assumptions:

$$\begin{aligned} E(\mathbf{f}) &= \mathbf{0}, \quad \text{Cov}(e_i, e_j) = 0 \quad (i \neq j), \\ E(\mathbf{e}) &= \mathbf{0}, \quad \text{Cov}(\mathbf{f}, \mathbf{e}) = \mathbf{0}, \\ \text{Var}(f_i) &= 1. \end{aligned}$$

With these assumptions, we can see that the covariance matrix of the error term has a diagonal form:

$$\text{Cov}(\mathbf{e}) = \mathbf{\Psi} = \text{Diag}(\psi_{11}, \dots, \psi_{pp}). \quad (6.2)$$

We can thus re-express the correlation matrix $\mathbf{\rho} = [(\rho_{ij})]$ of our random variables \mathbf{x} by our model:

$$\begin{aligned} \mathbf{\rho} &= \text{Cor}(\mathbf{x}) = \text{Cov}(\mathbf{y}) = \text{Cov}(\mathbf{\Lambda} \mathbf{f} + \mathbf{e}) \\ &= \mathbf{\Lambda} \underbrace{\text{Cov}(\mathbf{f})}_{\mathbf{\Phi}} \mathbf{\Lambda}^\top + \mathbf{\Lambda} \underbrace{\text{Cov}(\mathbf{f}, \mathbf{e})}_{\mathbf{0}} + \underbrace{\text{Cov}(\mathbf{e}, \mathbf{f})}_{\mathbf{0}} \mathbf{\Lambda}^\top + \underbrace{\text{Cov}(\mathbf{e})}_{\mathbf{\Psi}} \\ &= \mathbf{\Lambda} \mathbf{\Phi} \mathbf{\Lambda}^\top + \mathbf{\Psi}. \end{aligned}$$

Here, $\mathbf{\Phi}$ is the $(k \times k)$ matrix with the correlations between the factors. If we additionally assume that the factors are uncorrelated, i.e.

$$\text{Cov}(\mathbf{f}) = \mathbf{\Phi} = \mathbf{I}, \quad (6.3)$$

we obtain

$$\boldsymbol{\rho} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi} \quad (6.4)$$

or

$$\boldsymbol{\rho}_{red} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top = \begin{pmatrix} \kappa_1^2 & \rho_{12} & \cdots & \rho_{1p} \\ \rho_{21} & \kappa_2^2 & \cdots & \rho_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{p1} & \rho_{p2} & \cdots & \kappa_p^2 \end{pmatrix} = \boldsymbol{\rho} - \boldsymbol{\Psi} \quad , \quad (6.5)$$

with the *reduced correlation matrix* $\boldsymbol{\rho}_{red}$. The diagonal elements $\kappa_i^2 = 1 - \psi_{ii} = \sum_{j=1}^k \lambda_{ij}^2$ ($i = 1, \dots, p$) are called *communalities*. They correspond to the row-sums of the squared factor loadings, and describe the proportion of variance of y_i , explained by the factors, because the total variance is $\text{trace}(\boldsymbol{\rho})$, and thus $\text{trace}(\boldsymbol{\rho}_{red})$ is the variance explained by the factor model.

28 What are the uniquenesses and communalities, how are they defined, what is their meaning and how are they related?

6.2.2 Non-uniqueness of the factor loadings

Let \boldsymbol{G} be an orthogonal matrix of dimension $(k \times k)$. Because of $\boldsymbol{G}^{-1} = \boldsymbol{G}^\top$ we have

$$\boldsymbol{y} = (\boldsymbol{\Lambda}\boldsymbol{G})(\boldsymbol{G}^\top \boldsymbol{f}) + \boldsymbol{e} \quad . \quad (6.6)$$

Since the new factors $\boldsymbol{G}^\top \boldsymbol{f}$ also fulfill our model assumptions,

$$\begin{aligned} E(\boldsymbol{G}^\top \boldsymbol{f}) &= \mathbf{0}, & \text{Cov}(\boldsymbol{G}^\top \boldsymbol{f}) &= \boldsymbol{I} \\ \text{and} & & \text{Cov}(\boldsymbol{G}^\top \boldsymbol{f}, \boldsymbol{e}) &= \mathbf{0} \end{aligned}$$

this k -factor model is also valid with the new factors, i.e.

$$\boldsymbol{\rho} = (\boldsymbol{\Lambda}\boldsymbol{G})(\boldsymbol{G}^\top \boldsymbol{\Lambda}^\top) + \boldsymbol{\Psi} \quad . \quad (6.7)$$

This, however, means that the new factor loadings are not uniquely determined.

Basically, this is not a big problem, because later on we want to rotate the factors in any case in order to obtain a better interpretation. For now, if we just want to obtain a unique solution, we can impose further restrictions in order to achieve uniqueness. These are the constraints that either $\boldsymbol{\Lambda}^\top \boldsymbol{\Psi}^{-1} \boldsymbol{\Lambda}$ or $\boldsymbol{\Lambda}^\top \boldsymbol{\Lambda}$ is diagonal.

29 What is the maximum number of factors we can include in the factor model, and why?

6.2.3 Number of parameters

Practically, we can estimate the correlation matrix $\hat{\boldsymbol{\rho}}$ based on real data. Then, $\hat{\boldsymbol{\rho}}$ is used to estimate loadings $\hat{\boldsymbol{\Lambda}}$ and error variances $\hat{\boldsymbol{\Psi}}$, and they need to fulfill either

$$\hat{\boldsymbol{\Lambda}}^\top \hat{\boldsymbol{\Psi}}^{-1} \hat{\boldsymbol{\Lambda}} = \text{Diag} \quad \text{or} \quad \hat{\boldsymbol{\Lambda}}^\top \hat{\boldsymbol{\Lambda}} = \text{Diag} \quad (6.8)$$

as well as

$$\hat{\boldsymbol{\rho}} = \hat{\boldsymbol{\Lambda}} \hat{\boldsymbol{\Lambda}}^\top + \hat{\boldsymbol{\Psi}} \quad . \quad (6.9)$$

Does the factor model lead to a simpler interpretation than the correlation matrix? In other words, is the number of parameters for the factor model lower compared to the correlation matrix? For the correlation matrix we need to estimate $\frac{1}{2}p(p+1)$ parameters. For a k -factor model we need to estimate $\boldsymbol{\Lambda}$ and $\boldsymbol{\Psi}$, thus $pk + p$ parameters. For uniqueness of these estimates we ask for diagonal structure of either $\boldsymbol{\Lambda}^\top \boldsymbol{\Psi}^{-1} \boldsymbol{\Lambda}$ or $\boldsymbol{\Lambda}^\top \boldsymbol{\Lambda}$, thus $\frac{1}{2}k(k-1)$ restrictions. Therefore, for the k -factor model we have $pk + p - \frac{1}{2}k(k-1)$ parameters to estimate. The difference to an unrestricted model is thus

$$\begin{aligned} s &= \frac{1}{2}p(p+1) - \left(pk + p - \frac{1}{2}k(k-1) \right) \\ &= \frac{1}{2}(p-k)^2 - \frac{1}{2}(p+k) \quad . \end{aligned} \quad (6.10)$$

Only for the case $s > 0$, the k -factor model leads to a simpler interpretation than the correlation matrix, the other cases would not be valid (or useful) solutions. Thus, the requirement for $s > 0$ leads to an upper bound for the number k of factors.

How can we estimate the communalities and loadings (PFA)?

6.2.4 Parameter estimation by Principal Factor Analysis (PFA)

The main job in factor analysis is to estimate the parameters $\mathbf{\Lambda}$ and $\mathbf{\Psi}$. There are different approaches, such as the *Maximum Likelihood Method (MLM)*, implemented in the R function `factanal()`, or *Principal Factor Analysis (PFA)*, implemented in the R package `StatDA` as function `pfa()`. For MLM we need to assume that the data are generated from a multivariate normal distribution, which is not explicitly required for PFA. In fact, MLM is quite technical, and thus we focus here on PFA only, since this is closely related to PCA.

We start from our model

$$\boldsymbol{\rho} = \mathbf{\Lambda}\mathbf{\Lambda}^\top + \mathbf{\Psi} \quad (6.11)$$

and try to first estimate the communalities, which is equivalent to estimating the diagonal elements of $\mathbf{\Psi}$. Afterward we estimate the loadings matrix.

Estimation of the communalities

The communalities

$$\kappa_i^2 = \sum_{j=1}^k \lambda_{ij}^2 = 1 - \psi_{ii} \quad (i = 1, \dots, p) \quad (6.12)$$

describe the porportion of variance explained by the k -factor model. They are in the interval $[0, 1]$. There are different options to estimate the communalities:

1. Highest correlation coefficient: $\max_{i \neq j} |\hat{\rho}_{ij}|$

For the estimation of the communalities we use from each column of $\boldsymbol{\rho}$ the largest (absolute) non-diagonal element.

2. Squared multiple correlation coefficient: $\hat{\rho}_{i,12\dots i(\dots p)}^2$

This measure refers to an R^2 measure from a linear regression of the i -th variable on the remaining variables. Thus, it tells us about the variance porportion of the i -th variable explained from the other variables. We can compute this measure as:

2. Squared multiple correlation coefficient: $\hat{\rho}_{i,12\dots i(\dots p)}^2$

This measure refers to an R^2 measure from a linear regression of the i -th variable on the remaining variables. Thus, it tells us about the variance porportion of the i -th variable explained from the other variables. We can compute this measure as:

$$\hat{\rho}_{i,12\dots i(\dots p)}^2 = 1 - \frac{1}{\hat{\rho}^{ii}} \quad , \quad (6.13)$$

where $\hat{\rho}^{ii}$ is the i -th diagonal element of $\hat{\boldsymbol{\rho}}^{-1}$.

3. Iterative estimation:

We first need to fix the number k of factors. This could be done by using a criterion from PCA on the number of components. Then we start to initialize the communalities by using one of the estimation methods mentioned above. The diagonal elements of $\hat{\boldsymbol{\rho}}$ are replaced by these communalities, and from the resulting reduced correlation matrix we estimate the loadings (see end of this section). Based on that we re-estimate the communalities:

$$\hat{\kappa}_i^2 = \sum_{j=1}^k \hat{\lambda}_{ij}^2 \quad \text{for } i = 1, \dots, p \quad , \quad (6.14)$$

and they are used to form the new reduced correlation matrix. The iteration continues until the communalities stabilize.

If the communalities are over-estimated, part of the uniquenesses are forced into the k -factor model, which might change the pattern and interpretability of the factors. The same may happen if the communalities are under-estimated, because variance of the factors is forced into the uniquenesses. This issue is more severe if the number of variables is small.

Estimation of the loadings

Once the communalities $\hat{\kappa}_i^2 = 1 - \hat{\psi}_{ii}$, for $i = 1, \dots, p$, have been estimated, we can estimate the loadings based on the reduced correlation matrix

$$\hat{\rho}_{red} = \hat{\rho} - \hat{\Psi}, \quad (6.15)$$

where the diagonal consists of these communalities.

Using the Spectral Theorem 1.4.3, we have:

$$\hat{\rho}_{red} = \hat{\Gamma} \hat{\mathbf{A}} \hat{\Gamma}^\top, \quad (6.16)$$

where $\hat{\Gamma} = (\hat{\gamma}_1, \dots, \hat{\gamma}_p)$ and $\hat{\mathbf{A}} = \text{Diag}(\hat{a}_1, \dots, \hat{a}_p)$, with the eigenvalues \hat{a}_i and the corresponding eigenvectors $\hat{\gamma}_i$ of $\hat{\rho}_{red}$ ($i = 1, \dots, p$).

Note that we fixed the number of factors with k , where $1 \leq k < p$, and thus the loadings matrix needs to have dimension $p \times k$. Thus, we shall only use the first k eigenvectors $\hat{\Gamma}_{1:k}$ and eigenvalues $\hat{\mathbf{A}}_{1:k}$ in the k -factor model

$$\hat{\rho} - \hat{\Psi} = \hat{\Lambda} \hat{\Lambda}^\top = \hat{\Gamma}_{1:k} \hat{\mathbf{A}}_{1:k} \hat{\Gamma}_{1:k}^\top + \sum_{i=k+1}^p \hat{a}_i \hat{\gamma}_i \hat{\gamma}_i^\top. \quad (6.17)$$

The estimated loadings matrix is thus naturally

$$\hat{\Lambda} = \hat{\Gamma}_{1:k} \hat{\mathbf{A}}_{1:k}^{1/2}, \quad (6.18)$$

where the diagonal elements of $\hat{\mathbf{A}}_{1:k}^{1/2}$ are the values $\sqrt{\hat{a}_1}, \dots, \sqrt{\hat{a}_k}$.

Considering Equation (6.17), it is also natural to update the estimated uniquenesses as

$$\hat{\psi}_{ii} = 1 - \sum_{j=1}^k \hat{\lambda}_{ij}^2 \quad \text{for } i = 1, \dots, p. \quad (6.19)$$

The solution is valid if all $\hat{\psi}_{ii} \geq 0$.

31 How can we interpret factors? Give an overview of factor rotation criteria.

A rotation of the factors will change the loadings, and thus also the interpretation. The goal is to rotate in such a way that the resulting pattern of the loadings matrix is “simple” and thus interpretable. “Simple” basically means that the loadings matrix contains essentially small (absolute) values, and few values close to -1 or 1 . If there is a large (absolute) value, then we know that the corresponding variable has a strong contribution on this factor, while others with loadings close to zero do not have a contribution. Many loadings close to zero would simplify the interpretation of a factor. From the rotated factors we get strong contributions from some variables and weak contributions from others. The rotated factors would thus have a much clearer interpretation.

6.3.1 Orthogonal rotation

The interpretation is “simple” if a point in two dimensions is close to an axis. In that case, the product of both coordinates is small, and by taking the square one gets rid of the sign. Now we can consider the sum of all squared products as a criterion for simplicity. This should be valid for all pairs of factors, resulting in the criterion

$$\sum_{s < j=1}^k \sum_{i=1}^p (\lambda_{is} \lambda_{ij})^2 \longrightarrow \min \quad . \quad (6.20)$$

A factor rotation can be achieved by a transformation of the loadings matrix. For orthogonal rotation we thus have to consider an orthogonal $k \times k$ matrix \mathbf{T} , and the rotated loadings are given by

$$\tilde{\mathbf{\Lambda}} = \mathbf{\Lambda} \mathbf{T} \quad . \quad (6.21)$$

Orthogonal transformations do not change the communalities, since

$$\text{Diag}(\tilde{\mathbf{\Lambda}} \tilde{\mathbf{\Lambda}}^\top) = \text{Diag}(\mathbf{\Lambda} \mathbf{T} \mathbf{T}^{-1} \mathbf{\Lambda}^\top) = \text{Diag}(\mathbf{\Lambda} \mathbf{\Lambda}^\top) \quad ,$$

and thus

$$\kappa_i^2 = \sum_{j=1}^k \lambda_{ij}^2 = \sum_{j=1}^k \tilde{\lambda}_{ij}^2 \quad \text{for } i = 1, \dots, p \quad . \quad (6.22)$$

Also the squared communalities remain constant under orthogonal transformations,

$$(\kappa_i^2)^2 = \left(\sum_{j=1}^k \tilde{\lambda}_{ij}^2 \right)^2 = \sum_{j=1}^k \tilde{\lambda}_{ij}^4 + 2 \sum_{s < j=1}^k \tilde{\lambda}_{is}^2 \tilde{\lambda}_{ij}^2 = \text{const} \quad , \quad (6.23)$$

and also the sum over all variables remains constant:

$$\sum_{i=1}^p \sum_{j=1}^k \tilde{\lambda}_{ij}^4 + 2 \sum_{i=1}^p \sum_{s < j=1}^k \tilde{\lambda}_{is}^2 \tilde{\lambda}_{ij}^2 = \text{const} \quad . \quad (6.24)$$

Since the sum of both expressions is constant, one term is maximized if the other term is minimized, and vice versa.

The **Quartimax criterion** thus maximizes

$$QMAX = \sum_{i=1}^p \sum_{j=1}^k \tilde{\lambda}_{ij}^4, \quad (6.25)$$

with the effect that this may lead to one dominant factor.

Another criterion for factor rotation is the **Varimax criterion** (Kaiser, 1958), where the variance of the squared factor loadings for each factor j is considered:

$$s_j^2 = \frac{1}{p} \sum_{i=1}^p \left(\tilde{\lambda}_{ij}^2 - \frac{1}{p} \sum_{l=1}^p \tilde{\lambda}_{lj}^2 \right)^2 = \frac{1}{p} \sum_{i=1}^p (\tilde{\lambda}_{ij}^2)^2 - \frac{1}{p^2} \left[\sum_{i=1}^p \tilde{\lambda}_{ij}^2 \right]^2 \quad (6.26)$$

This variance is summed up over all factors. Maximizing this expression means that we want to have (absolute) large and small loadings. Since variables with larger communality are dominating the criterion, one can normalize with the communalities, which yields (after multiplication with p^2) the varimax criterion

$$VMAX = p \sum_{j=1}^k \sum_{i=1}^p \left(\frac{\tilde{\lambda}_{ij}}{\kappa_i} \right)^4 - \sum_{j=1}^k \left[\sum_{i=1}^p \left(\frac{\tilde{\lambda}_{ij}}{\kappa_i} \right)^2 \right]^2 \quad . \quad (6.27)$$

6.3.2 Oblique rotation

The main difference to orthogonal rotations is in the initial factor model, which is now

$$\boldsymbol{\rho} = \boldsymbol{\Lambda} \boldsymbol{\Phi} \boldsymbol{\Lambda}^\top + \boldsymbol{\Psi} \quad , \quad (6.28)$$

because oblique factors are no longer uncorrelated, and thus the correlation matrix $\boldsymbol{\Phi}$ of the factors needs to be considered.

The **Quartimin criterion** is defined as in Equation (6.20) by

$$QMIN = \sum_{s < j=1}^k \sum_{i=1}^p \tilde{\lambda}_{is}^2 \tilde{\lambda}_{ij}^2 \quad . \quad (6.29)$$

The rotated loadings are

$$\tilde{\boldsymbol{\Lambda}} = \boldsymbol{\Lambda} \boldsymbol{T} \quad , \quad (6.30)$$

where \boldsymbol{T} is no longer orthogonal. Plugging in the rotated loadings in the model yields

$$\begin{aligned} \boldsymbol{\rho} - \boldsymbol{\Psi} &= \boldsymbol{\Lambda} \boldsymbol{\Phi} \boldsymbol{\Lambda}^\top = \tilde{\boldsymbol{\Lambda}} \boldsymbol{T}^{-1} \boldsymbol{\Phi} (\boldsymbol{T}^{-1})^\top \tilde{\boldsymbol{\Lambda}}^\top = \tilde{\boldsymbol{\Lambda}} \boldsymbol{T}^{-1} \text{Cov}(\boldsymbol{f}) (\boldsymbol{T}^{-1})^\top \tilde{\boldsymbol{\Lambda}}^\top \\ &= \tilde{\boldsymbol{\Lambda}} \text{Cov}(\boldsymbol{T}^{-1} \boldsymbol{f}) \tilde{\boldsymbol{\Lambda}}^\top = \tilde{\boldsymbol{\Lambda}} \text{Cov}(\tilde{\boldsymbol{f}}) \tilde{\boldsymbol{\Lambda}}^\top \quad . \end{aligned} \quad (6.31)$$

Die rotated factors $\tilde{\boldsymbol{f}}$ are thus

$$\tilde{\boldsymbol{f}} = \boldsymbol{T}^{-1} \boldsymbol{f} \quad . \quad (6.32)$$

The transformation matrix \boldsymbol{T} thus needs to be invertible.

A more flexible criterion is the **Oblimin criterion**

$$OBMIN = \sum_{s < j=1}^k \left(\sum_{i=1}^p \tilde{\lambda}_{is}^2 \tilde{\lambda}_{ij}^2 - \frac{\gamma}{p} \sum_{i=1}^p \tilde{\lambda}_{is}^2 \sum_{i=1}^p \tilde{\lambda}_{ij}^2 \right) \quad . \quad (6.33)$$

The choice $\gamma = 0$ yields the Quartimin criterion, while $\gamma = 1$ leads to the so-called *Covarimin criterion*.

32 How are factor scores estimated (Bartlett and Regression method). Name the models, formulas and solutions for the factor scores estimates.

6.4.1 Weighted least-squares estimation

In the R function `factanal()` this method is selected by `scores="Bartlett"`.

The factor model is

$$\boldsymbol{y} = \boldsymbol{\Lambda} \boldsymbol{f} + \boldsymbol{e} \quad (6.34)$$

with the error term $\boldsymbol{e} = (e_1, \dots, e_p)^\top$. This model could also be considered as a regression model, by regressing \boldsymbol{y} on $\boldsymbol{\Lambda}$, with the regression coefficients \boldsymbol{f} . However, the model is heteroscedastic, since the error variances $\text{Var}(e_i) = \psi_i$ for $i = 1, \dots, p$ are not necessarily equal. However, one can multiply the equation with weights $\boldsymbol{\Psi}^{-1/2}$, which yields

$$\boldsymbol{\Psi}^{-1/2} \boldsymbol{y} = \boldsymbol{\Psi}^{-1/2} \boldsymbol{\Lambda} \boldsymbol{f} + \boldsymbol{\Psi}^{-1/2} \boldsymbol{e} \quad . \quad (6.35)$$

The covariance of the new error term is the identity matrix, and thus we have a homoscedastic model, and the least-squares estimator can be used, where the loadings and uniquenesses are considered as “true” values. This results in

$$\hat{\boldsymbol{f}} = (\boldsymbol{\Lambda}^\top \boldsymbol{\Psi}^{-1} \boldsymbol{\Lambda})^{-1} \boldsymbol{\Lambda}^\top \boldsymbol{\Psi}^{-1} \boldsymbol{y} \quad (6.36)$$

for the estimated factors (random variables).

For a given $n \times p$ data matrix \mathbf{X} we obtain the mean-centered and scaled matrix \mathbf{Y} as a realization of \mathbf{y} , which can be substituted into Equation (6.36), to obtain the $n \times k$ matrix of estimated factor scores:

$$\hat{\mathbf{F}} = \mathbf{Y}\Psi^{-1}\Lambda (\Lambda^\top \Psi^{-1}\Lambda)^{-1} \quad (6.37)$$

6.4.2 Regression method

In the R function `factanal()` this method is selected by `scores="regression"`.

As for the previous method, the loadings Λ and uniquenesses Ψ are considered as known. Here we again consider a regression problem, but we regress the (unknown) factors \mathbf{f} on \mathbf{y} (multivariate regression):

$$\mathbf{f} = \mathbf{B}\mathbf{y} + \boldsymbol{\delta} , \quad (6.38)$$

where \mathbf{B} is the $k \times p$ matrix of regression coefficients, and $\boldsymbol{\delta}$ the error matrix. The least-squares estimator is

$$\hat{\mathbf{B}} = \mathbf{f}\mathbf{y}^\top (\mathbf{y}\mathbf{y}^\top)^{-1} , \quad (6.39)$$

and the estimated factors are

$$\hat{\mathbf{f}} = \hat{\mathbf{B}}\mathbf{y} = \mathbf{f}\mathbf{y}^\top (\mathbf{y}\mathbf{y}^\top)^{-1}\mathbf{y} . \quad (6.40)$$

Now we have the factors \mathbf{f} again on the right-hand side, but we can plug in our factor model for \mathbf{y} :

$$\begin{aligned} \hat{\mathbf{f}} &= \mathbf{f}(\Lambda\mathbf{f} + \mathbf{e})^\top (\mathbf{y}\mathbf{y}^\top)^{-1}\mathbf{y} \\ &= (\mathbf{f}\mathbf{f}^\top \Lambda^\top + \mathbf{f}\mathbf{e}^\top)(\mathbf{y}\mathbf{y}^\top)^{-1}\mathbf{y} \\ &= \mathbf{f}\mathbf{f}^\top \Lambda^\top (\mathbf{y}\mathbf{y}^\top)^{-1}\mathbf{y} \end{aligned}$$

Substituting the sample version yields

$$\begin{aligned} \hat{\mathbf{F}} &= \mathbf{Y}(n-1)(\mathbf{Y}^\top \mathbf{Y})^{-1}\Lambda \frac{1}{n-1}\mathbf{F}^\top \mathbf{F} \\ &= \mathbf{Y}\mathbf{R}^{-1}\Lambda \hat{\Phi} , \end{aligned}$$

where $\hat{\Phi}$ is the estimated correlation matrix of the factors. In the orthogonal case, the factors are uncorrelated, and thus

$$\hat{\mathbf{F}} = \mathbf{Y}\mathbf{R}^{-1}\Lambda . \quad (6.41)$$

33 What is the problem setting in multiple correlation analysis? What is the objective function to minimize?

In statistical data analysis, one is often interested in determining relationships and dependencies of features. If one also wants to measure the existence and the strength of dependencies, correlation analysis can be used. The correlation measures the linear relationship between features.

7.1 Multiple correlation analysis

The *multiple correlation* is a measure of the dependency of a feature x on a p - dimensional feature $\mathbf{y} = (y_1, \dots, y_p)^\top$. We assume that both x and y_1, \dots, y_p are random variables with a joint distribution. The mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ of this distribution (not necessarily normal distribution) are

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_x \\ \boldsymbol{\mu}_y \end{pmatrix} \quad \text{bzw.} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{xx} & \boldsymbol{\sigma}_{y_x}^\top \\ \boldsymbol{\sigma}_{y_x} & \boldsymbol{\Sigma}_{yy} \end{pmatrix},$$

where

$$\begin{aligned} E(x) &= \mu_x & \text{Var}(x) &= \sigma_{xx} \\ E(\mathbf{y}) &= \boldsymbol{\mu}_y & \text{Cov}(\mathbf{y}) &= \boldsymbol{\Sigma}_{yy} \end{aligned}$$

and

$$\text{Cov}(\mathbf{y}, x) = \boldsymbol{\sigma}_{y_x} \quad \text{Cov}(x, \mathbf{y}) = \boldsymbol{\sigma}_{y_x}^\top.$$

If x is now predicted by \mathbf{y} (linear), the error is analogous to regression analysis

$$x - a_0 - a_1 y_1 - \dots - a_p y_p.$$

This error is random, and therefore one would like to choose the coefficients a_0 and $\mathbf{a} = (a_1, \dots, a_p)^\top$ such that the *mean squared error* (MSE)

$$\text{MSE} = E(x - a_0 - \mathbf{a}^\top \mathbf{y})^2$$

is minimal. However, this MSE depends on the joint distribution of x and \mathbf{y} , and thus on the parameters $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$.

- 34 What is the linear prediction function in multiple correlation analysis? Describe the structure of the proof.**

Theorem 7.1.1 The linear prediction function $a_0 + \mathbf{a}^\top \mathbf{y}$ with the coefficients

$$\mathbf{a} = \Sigma_{\mathbf{y}\mathbf{y}}^{-1} \sigma_{\mathbf{y}x} \quad \text{and} \quad a_0 = \mu_x - \mathbf{a}^\top \mu_{\mathbf{y}}$$

has minimal MSE of all linear prediction functions of x . It holds that

$$MSE = E(x - a_0 - \mathbf{a}^\top \mathbf{y})^2 = E(x - \mu_x - \sigma_{\mathbf{y}x}^\top \Sigma_{\mathbf{y}\mathbf{y}}^{-1} (\mathbf{y} - \mu_{\mathbf{y}}))^2 = \sigma_{xx} - \sigma_{\mathbf{y}x}^\top \Sigma_{\mathbf{y}\mathbf{y}}^{-1} \sigma_{\mathbf{y}x} .$$

Furthermore,

$$a_0 + \mathbf{a}^\top \mathbf{y} = \mu_x + \sigma_{\mathbf{y}x}^\top \Sigma_{\mathbf{y}\mathbf{y}}^{-1} (\mathbf{y} - \mu_{\mathbf{y}})$$

is the linear prediction function that has maximum correlation with x , namely

$$\text{Corr}(x, a_0 + \mathbf{a}^\top \mathbf{y}) = \sqrt{\frac{\mathbf{a}^\top \Sigma_{\mathbf{y}\mathbf{y}} \mathbf{a}}{\sigma_{xx}}} = \sqrt{\frac{\sigma_{\mathbf{y}x}^\top \Sigma_{\mathbf{y}\mathbf{y}}^{-1} \sigma_{\mathbf{y}x}}{\sigma_{xx}}} .$$

Proof: See, for example, Johnson and Wichern (1998). For the proof of the minimal MSE one needs to re-express the squared expectation, and then it is visible that with these choices of \mathbf{a} and a_0 one obtains a minimum of the MSE. For the maximum correlation one can make use of the

Extended Cauchy-Schwarz Inequality: Let \mathbf{b} and \mathbf{d} be two vectors, and \mathbf{B} a positive definite matrix. Then

$$(\mathbf{b}^\top \mathbf{d})^2 \leq (\mathbf{b}^\top \mathbf{B} \mathbf{b})(\mathbf{d}^\top \mathbf{B}^{-1} \mathbf{d}) , \quad (7.1)$$

with equality if and only if $\mathbf{b} = c \mathbf{B}^{-1} \mathbf{d}$ (or $\mathbf{d} = c \mathbf{B} \mathbf{b}$), for a constant c .

The correlation between x and the best linear prediction function is called *multiple correlation coefficient* (of the population)

$$\rho_{x,\mathbf{y}} = + \sqrt{\frac{\sigma_{\mathbf{y}x}^\top \Sigma_{\mathbf{y}\mathbf{y}}^{-1} \sigma_{\mathbf{y}x}}{\sigma_{xx}}} .$$

Its squared form $\rho_{x,\mathbf{y}}^2$ is called *multiple coefficient of determination* (of the population), and it indicates how well the feature x is explained by the properties $\mathbf{y} = (y_1, \dots, y_p)^\top$.

If correlations are given instead of the covariances, then the multiple correlation coefficient can also be defined as

$$\rho_{x,\mathbf{y}}^2 = \rho_{\mathbf{y}x}^\top \rho_{\mathbf{y}\mathbf{y}}^{-1} \rho_{\mathbf{y}x} .$$

For a specific sample, one can write the theoretical quantities by the corresponding realizations, and hence obtain for the multiple correlation coefficient

$$r_{x,\mathbf{y}}^2 = \mathbf{R}_{\mathbf{y}x}^\top \mathbf{R}_{\mathbf{y}\mathbf{y}}^{-1} \mathbf{R}_{\mathbf{y}x} .$$

35 Name a hypothesis test for the multiple correlation coefficient.

A test for the hypothesis that the multiple correlation is zero is equivalent to making all bivariate correlations zero. So, we test the hypothesis

$$H_0 : \rho_{x,\mathbf{y}} = 0 \quad (= \rho_{xy_1} = \dots = \rho_{xy_p})$$

against the alternative hypothesis

$$H_1 : \exists i \in \{1, \dots, p\} \quad \text{with} \quad \rho_{xy_i} \neq 0$$

at the significance level α . If we can assume multivariate normal distribution, and n is the number of observations in the sample, the test statistic

$$F = \frac{(n-1-p) r_{x,y}^2}{p (1 - r_{x,y}^2)}$$

has an $F_{p,n-1-p}$ -distribution and the null hypothesis is rejected at the significance level of α , if

$$F > F_{p,n-1-p;1-\alpha}.$$

36 What is the problem setting in canonical correlation, what is the maximization problem?

In canonical correlation analysis we are interested in the linear dependence between two groups of variables. As it turns out, this dependence can no longer be expressed by a single correlation coefficient, but results in a subspace that describes the linear dependence between the groups.

Theorem 7.2.1 *Let \mathbf{x} be a p -dimensional and \mathbf{y} a q -dimensional random variable ($p \leq q$) with the expected values*

$$E(\mathbf{x}) = \boldsymbol{\mu}_1 \quad \text{and} \quad E(\mathbf{y}) = \boldsymbol{\mu}_2 .$$

The covariance matrices $\boldsymbol{\Sigma}_{ij}$ with $i, j = 1, 2$ are defined by

$$\begin{aligned} \boldsymbol{\Sigma}_{11} &= E[(\mathbf{x} - \boldsymbol{\mu}_1)(\mathbf{x} - \boldsymbol{\mu}_1)^\top] \quad , \\ \boldsymbol{\Sigma}_{22} &= E[(\mathbf{y} - \boldsymbol{\mu}_2)(\mathbf{y} - \boldsymbol{\mu}_2)^\top] \quad \text{and} \\ \boldsymbol{\Sigma}_{12} &= \boldsymbol{\Sigma}_{21}^\top = E[(\mathbf{x} - \boldsymbol{\mu}_1)(\mathbf{y} - \boldsymbol{\mu}_2)^\top] \end{aligned}$$

and have full rank. We consider the linear combinations $\varphi = \mathbf{a}^\top \mathbf{x}$ and $\eta = \mathbf{b}^\top \mathbf{y}$, where \mathbf{a} is a p -dimensional and \mathbf{b} is a q -dimensional vector.

Then the simple correlation between φ and η is given by

$$\max_{\mathbf{a}, \mathbf{b}} \text{Corr}(\varphi, \eta) = \rho_1 .$$

37 How do we get the linear combinations for canonical correlation? Why is a matrix product and Eigenvector/Eigenvalue problem involved?

The maximum is achieved by the linear combinations

$$\varphi_1 = \underbrace{\mathbf{e}_1^\top \Sigma_{11}^{-1/2}}_{\mathbf{a}_1^\top} \mathbf{x} \quad \text{and} \quad \eta_1 = \underbrace{\mathbf{f}_1^\top \Sigma_{22}^{-1/2}}_{\mathbf{b}_1^\top} \mathbf{y},$$

which are referred to as the first pair of canonical variables. ρ_1 is called the first canonical correlation coefficient.

The k -th pair of canonical variables ($k = 2, 3, \dots, p$) is given by

$$\varphi_k = \underbrace{\mathbf{e}_k^\top \Sigma_{11}^{-1/2}}_{\mathbf{a}_k^\top} \mathbf{x} \quad \text{and} \quad \eta_k = \underbrace{\mathbf{f}_k^\top \Sigma_{22}^{-1/2}}_{\mathbf{b}_k^\top} \mathbf{y}$$

and maximize

$$\text{Corr}(\varphi_k, \eta_k) = \rho_k$$

over all linear combinations that are uncorrelated with the previous $1, 2, \dots, k-1$ canonical variables. ρ_k is called k -th canonical correlation coefficient.

$\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_p^2$ are eigenvalues of the matrix $\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2}$ and $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p$ are the respective eigenvectors (dimension p).

$\rho_1^2, \rho_2^2, \dots, \rho_p^2$ are the p largest eigenvalues of the matrix $\Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1/2}$ with the respective eigenvectors $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_p$ (dimension q). Every \mathbf{f}_i is proportional to $\Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1/2} \mathbf{e}_i$.

The canonical variables have the following properties

$$\begin{aligned} \text{Var}(\varphi_k) &= \text{Var}(\eta_k) = 1 \\ \text{Cov}(\varphi_k, \varphi_l) &= \text{Corr}(\varphi_k, \varphi_l) = 0 & k \neq l \\ \text{Cov}(\eta_k, \eta_l) &= \text{Corr}(\eta_k, \eta_l) = 0 & k \neq l \\ \text{Cov}(\varphi_k, \eta_l) &= \text{Corr}(\varphi_k, \eta_l) = 0 & k \neq l \end{aligned}$$

for $k, l = 1, 2, \dots, p$; and with the notation $\boldsymbol{\varphi} = (\varphi_1, \dots, \varphi_p)^\top$, $\boldsymbol{\eta} = (\eta_1, \dots, \eta_p)^\top$ and $\boldsymbol{\rho} = \text{Diag}(\rho_1, \dots, \rho_p)$, it follows that

$$\text{Cov} \begin{pmatrix} \boldsymbol{\varphi} \\ \boldsymbol{\eta} \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \boldsymbol{\rho} \\ \boldsymbol{\rho} & \mathbf{I} \end{pmatrix}.$$

The canonical correlation coefficients are invariant to linear transformations, which will be shown below.

38 What happens if there is the same variable in X and Y in canonical correlation?

If the same variable is included in both sets XXX and YYY, the first canonical correlation will likely be close to 1, as that variable will have a perfect correlation with itself. This redundancy can lead to overestimation of the strength of the relationship between the two variable sets. It is typically advisable to remove duplicate variables to avoid misleading results in canonical correlation analysis.

Was bedeutet ein kanonischer Korrelationskoeffizient ρ_1 von 1? Es reicht bereits, dass 1 Variable x_i aus \mathbf{x} gleich einer Variable y_j aus \mathbf{y} ist; ist (oBdA?) z.B. $x_1 = y_1$, (und die übrigen Einträge in \mathbf{x} und \mathbf{y} beliebig?) kann man mit Vektoren $\mathbf{a}_1 = (1, 0, \dots, 0)^T$ und $\mathbf{b}_1 = (1, 0, \dots, 0)^T$ linear kombinieren, so dass sich ein $\rho_1 = 1$ ergibt. Man kann dann nichts über die anderen Variablen aus \mathbf{x} und \mathbf{y} aussagen. (Kommentar: habe das so in etwa bei der Prüfung erklärt bekommen).

**39 What are some hypothesis tests in canonical correlation analysis?
(WS20: focus on permutation test)**

In canonical correlation analysis, there are a number of tests, of which only selected ones are mentioned here, because the distributions of the estimators are very complicated. Assuming that the data are normally distributed, the following applies:

- (a) A likelihood ratio test for the hypothesis $H_0 : \Sigma_{12} = \mathbf{O}$, i.e. for the hypothesis that \mathbf{x} and \mathbf{y} are uncorrelated, is given by the test statistic

$$\lambda^{2/n} = |\mathbf{I} - \mathbf{S}_{22}^{-1} \mathbf{S}_{21} \mathbf{S}_{11}^{-1} \mathbf{S}_{12}| = \prod_{i=1}^p (1 - r_i^2).$$

This test statistic has a $\Lambda(q, n - 1 - p, p)$ - Wilks distribution. n is the sample size and r_1, \dots, r_p are the sample canonical correlation coefficients.

- (b) Since the Wilks distribution as a product of beta distributions is relatively complex, it can be approximated (with Bartlett's approximation) for large sample sizes n by

$$- \left[n - \frac{1}{2}(p + q + 3) \right] \ln \prod_{i=1}^p (1 - r_i^2) \sim \chi_{pq}^2.$$

- (c) A similar test statistic can also be formulated for the hypothesis that only s of the canonical correlation coefficients are non-zero, namely by

$$- \left[n - \frac{1}{2}(p + q + 3) \right] \ln \prod_{i=s+1}^p (1 - r_i^2) \sim \chi_{(p-s)(q-s)}^2.$$

A permutation test is carried out, which does not rely on a distributional assumption. The idea is to permute the observations of one data set, by keeping the other data set unpermuted. Then the canonical correlations are estimated, and this is done many times. Finally, a p-value is determined as the fraction of bootstrap correlation results exceeding the canonical correlation of the unpermuted data.

40 What is the goal of discriminant analysis? What is the expected cost of misclassification, what is involved? How can the ECM be minimized, and how do we arrive at those rules?

It is a multivariate method that deals on the one hand with the classification of different object groups and on the other hand with the assignment of new

objects to previously determined groups. In the former case, the attempt is made to capture the differences of the objects which are known to originate from two or more populations, either graphically or algebraically. One is thus looking for a discriminant function that allows the best possible separation. In the second case, one would like to divide the objects into two or more groups. The goal is then to classify new objects by means of defined rules. The above cases are often directly related, because a function that separates objects can also be used to classify new objects or vice versa.

The objects are classified on the basis of measurements of p underlying random variables $X = (X_1, \dots, X_p)^T$. Assuming that there are two groups, the objects to be measured are then divided into the classes π_1 and π_2 respectively. Let the sum of the values of the first class be the population of the x -values of π_1 , and that of the second class the population of the x -values of π_2 . The two populations are then described by the probability distributions $f_1(x)$ and $f_2(x)$.

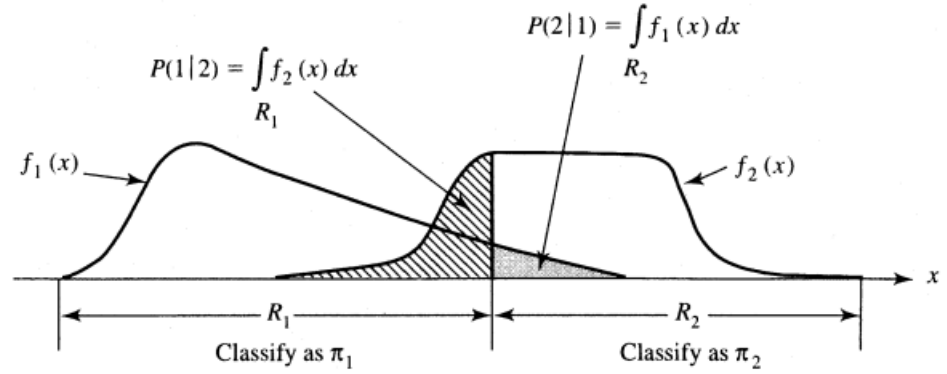
Let Ω be the entire sampling space, i.e. the space containing all observations x . Each object x must come from either population π_1 or π_2 . Furthermore, let R_1 be the space of observations x to which we assign the objects of π_1 , and R_2 the space to which the remaining objects of π_2 are assigned. Ω is the union of R_1 and R_2 . It may happen in the case of group assignment that objects which actually belong to population π_1 are falsely classified as π_2 . If the probability functions $f_1(x)$ and $f_2(x)$ are known, this probability of incorrect assignment can be calculated as a conditional probability $P(2|1)$ by:

$$P(2|1) = P(\mathbf{X} \in R_2|\pi_1) = \int_{R_2=\Omega-R_1} f_1(\mathbf{x})d\mathbf{x} . \quad (8.1)$$

Conversely, it may be the case that objects originating from population π_2 are erroneously assigned to π_1 . The corresponding probability is

$$P(1|2) = P(\mathbf{X} \in R_1|\pi_2) = \int_{R_1} f_2(\mathbf{x})d\mathbf{x} . \quad (8.2)$$

The integral in (8.1) describes the volume of the density function $f_1(\mathbf{x})$ over the region R_2 , or analogously for (8.2).



Let p_1 be the probability that the objects come from π_1 (*prior probability*), and p_2 those for π_2 , where $p_1 + p_2 = 1$ must hold. Then, the following probabilities can be calculated by applying the formula for conditional probabilities:

$$\begin{aligned} P(\text{Observation correctly classified as } \pi_1) &= P(\mathbf{X} \in R_1|\pi_1)P(\pi_1) = P(1|1)p_1 \\ P(\text{Observation incorrectly classified as } \pi_1) &= P(\mathbf{X} \in R_1|\pi_2)P(\pi_2) = P(1|2)p_2 \\ P(\text{Observation correctly classified as } \pi_2) &= P(\mathbf{X} \in R_2|\pi_2)P(\pi_2) = P(2|2)p_2 \\ P(\text{Observation incorrectly classified as } \pi_2) &= P(\mathbf{X} \in R_2|\pi_1)P(\pi_1) = P(2|1)p_1 \end{aligned}$$

Missclassification is often directly associated with costs. Of course, the costs are 0 if correctly classified. They are $c(1|2)$ if an observation of π_2 is wrongly classified as π_1 . And the cost is $c(2|1)$ if observations of π_1 are erroneously classed as π_2 . Together with the probabilities for misclassification, the expected costs for misclassification (ECM) can now be calculated as

$$ECM = c(2|1)P(2|1)p_1 + c(1|2)P(1|2)p_2 . \quad (8.3)$$

The goal of a classification rule is to keep ECM as small as possible.

Theorem 8.2.1 *A classification rule that minimizes ECM is as follows: The set R_1 is defined for observations \mathbf{x} , for which the following applies:*

$$\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \quad (8.4)$$

The set R_2 is defined for observations \mathbf{x} , for which the following applies:

$$\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \quad (8.5)$$

Theorem 8.2.1 leads to the following special cases:

(a) $p_1 = p_2$:

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{c(1|2)}{c(2|1)} \quad R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{c(1|2)}{c(2|1)}$$

(b) $c(1|2) = c(2|1)$:

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{p_2}{p_1} \quad R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{p_2}{p_1} \quad (8.6)$$

(c) $p_1 = p_2$ and $c(1|2) = c(2|1)$:

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq 1 \quad R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < 1$$

However, other criteria can also be used to create a classification rule. For example, one could choose R_1 and R_2 so that the total probability of misclassification (TPM) is minimal, i.e.

$$\begin{aligned} TPM &= P(\text{Misclassification of an observation from } \pi_1 \text{ or } \pi_2) \\ &= p_1 \int_{R_2} f_1(\mathbf{x}) d\mathbf{x} + p_2 \int_{R_1} f_2(\mathbf{x}) d\mathbf{x} \end{aligned} \quad (8.7)$$

$$= p_1 P(2|1) + p_2 P(1|2) . \quad (8.8)$$

Assuming that the cost of misclassification is the same, one immediately realizes that minimizing (8.8) is equivalent to minimizing (8.3).

41 Two-group case: What is linear discriminant analysis, what are the assumptions? How to arrive at the rule for classification? Why is it called linear discriminant analysis, and how do we estimate the involved components?

The rule in case of p-dimensional features reduces the decision to a 1-dimensional variable y, which results from the corresponding linear combinations of the observations of π_1 and π_2 .

42 Two-group case: Explain the downprojection in LDA (graphic would probably help), if the priors and costs are equal.

8.3 The two-group case

We limit ourselves here to multivariate normally distributed populations. Thus, $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ are density functions of multivariate normal distributions with means $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$ and covariance matrices $\boldsymbol{\Sigma}_1$ and $\boldsymbol{\Sigma}_2$, respectively.

8.3.1 The special case $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}$

The joint density function of the random variable $\mathbf{X} = (X_1, \dots, X_p)^\top$ for the populations π_1 and π_2 are given by

$$f_i(\mathbf{x}) = (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right\} \quad \text{for } i = 1, 2. \quad (8.9)$$

If the parameters $\boldsymbol{\mu}_1$, $\boldsymbol{\mu}_2$ and $\boldsymbol{\Sigma}$ are known, then according to (8.4), the region that minimizes ECM is:

$$\begin{aligned} R_1 : \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \right\} \\ \geq \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \\ R_2 : \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \right\} \\ < \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \end{aligned} \quad (8.10)$$

The following classification rule can be specified using these regions R_1 and R_2 :

Theorem 8.3.1 *Let π_1 and π_2 be normally distributed populations with the same covariance $\boldsymbol{\Sigma}$. Then the classification rule for minimizing ECM is:*

An observation \mathbf{x}_0 is assigned to π_1 , if

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1} \mathbf{x}_0 - \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2) \geq \ln \left(\frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1} \right). \quad (8.11)$$

Otherwise \mathbf{x}_0 is assigned to π_2 .

42. Two-group case: Downprojection in LDA when priors and costs are equal

- Linear Discriminant Analysis (LDA) projects data onto a one-dimensional axis to maximize class separability.
- If priors and misclassification costs are equal, the decision boundary is determined by the mean vectors and pooled covariance matrix.
- The downprojection means that each sample is projected onto the discriminant axis w , where:

$$w = \Sigma^{-1}(\mu_1 - \mu_2)$$

with Σ being the within-class covariance matrix and μ_1, μ_2 the class means.

- This projection maximizes the ratio of between-class variance to within-class variance, ensuring optimal classification.

43 Two-group case: Why is QDA called "quadratic"?

8.3.2 The special case $\Sigma_1 \neq \Sigma_2$

The previously formulated classification rules were based on (8.4) and looked at the relation of the density functions $f_1(\mathbf{x})/f_2(\mathbf{x})$. In the case of the same covariance matrices, this ratio is reduced to a relatively simple term, which is usually expressed by (8.11). In the case of unequal covariance matrices (and unequal means), however, the ratio of density functions becomes a more complicated expression. As with (8.11), the logarithm of the relation can also be considered here:

$$\begin{aligned} \ln \left(\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \right) &= \ln \left(\frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} \right) - \frac{1}{2}(\mathbf{x} - \mu_1)^\top \Sigma_1^{-1}(\mathbf{x} - \mu_1) + \frac{1}{2}(\mathbf{x} - \mu_2)^\top \Sigma_2^{-1}(\mathbf{x} - \mu_2) \\ &= \frac{1}{2} \ln \left(\frac{|\Sigma_2|}{|\Sigma_1|} \right) - \frac{1}{2} \mathbf{x}^\top (\Sigma_1^{-1} - \Sigma_2^{-1}) \mathbf{x} + (\mu_1^\top \Sigma_1^{-1} - \mu_2^\top \Sigma_2^{-1}) \mathbf{x} \end{aligned}$$

$$-\frac{1}{2}\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2$$

This results in the following classification rule:

$$\begin{aligned} R_1 : & -\frac{1}{2}\mathbf{x}^\top (\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1})\mathbf{x} + (\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}_2^{-1})\mathbf{x} - k \geq \ln \left(\frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1} \right) \\ R_2 : & -\frac{1}{2}\mathbf{x}^\top (\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1})\mathbf{x} + (\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}_2^{-1})\mathbf{x} - k < \ln \left(\frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1} \right) \end{aligned}$$

with

$$k = \frac{1}{2} \ln \left(\frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_2|} \right) + \frac{1}{2} (\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2) .$$

The constant k now only depends on the mean and covariance of the two distributions. In the case of $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2$, the above rule is reduced to (8.11). The following statement now follows directly:

Theorem 8.3.3 *Let π_1 and π_2 be normally distributed populations with mean vectors $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$ and covariance matrices $\boldsymbol{\Sigma}_1$ and $\boldsymbol{\Sigma}_2$, respectively. Then, the classification rule for minimizing ECM is:*

An observation \mathbf{x}_0 is assigned to π_1 , if

$$-\frac{1}{2}\mathbf{x}_0^\top (\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1})\mathbf{x}_0 + (\boldsymbol{\mu}_1^\top \boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\mu}_2^\top \boldsymbol{\Sigma}_2^{-1})\mathbf{x}_0 - k \geq \ln \left(\frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1} \right) . \quad (8.16)$$

Otherwise \mathbf{x}_0 is assigned to π_2 .

The above discriminant function is *quadratic* in \mathbf{x} . Since the population parameters $\boldsymbol{\mu}_1$, $\boldsymbol{\mu}_2$, $\boldsymbol{\Sigma}_1$ and $\boldsymbol{\Sigma}_2$ are not known, they need to be obtained by the respective estimators, e.g. by the empirical sample estimators $\bar{\mathbf{x}}_1$, $\bar{\mathbf{x}}_2$, \mathbf{S}_1 and \mathbf{S}_2

43. Why is QDA called "quadratic"?

- Quadratic Discriminant Analysis (QDA) extends LDA by allowing each class to have its own covariance matrix $\boldsymbol{\Sigma}_k$.
- The decision boundary is no longer linear because the classification rule is based on the log-determinant and inverse of $\boldsymbol{\Sigma}_k$, leading to a quadratic form:

$$\mathbf{x}^\top \boldsymbol{\Sigma}_k^{-1} \mathbf{x} + \text{linear terms} + \text{constant} = 0$$

- This results in curved decision boundaries, making QDA more flexible than LDA but also more prone to overfitting when data is limited.

44 Two-group case: What is the Fischer criterion to maximize? What is the solution for the projection vector? What is the relation to LDA?

Fisher (1938) developed a linear discriminant function analogously to (8.14), but the idea behind it was different. He tried to transform multivariate observations to univariate, so that the two transformed groups are as strongly separated as possible. The idea was thus to find a direction $\mathbf{a} \in \mathbb{R}^p$, and to project the observations $\mathbf{a}^\top \mathbf{x}$ in order to obtain univariate values $y_{11}, y_{12}, \dots, y_{1n_1}$ for the observations of the first group and values $y_{21}, y_{22}, \dots, y_{2n_2}$ for the observations of the second group. The separation of the two populations then occurs in such a way that the arithmetic means \bar{y}_1 and \bar{y}_2 of the univariate y values deviate as much as possible. This difference is expressed in units of the standard deviation and is therefore used as a criterion for the separation

$$\frac{|\bar{y}_1 - \bar{y}_2|}{s_y} \longrightarrow \max$$

with the pooled variance of the y -values

$$s_y^2 = \frac{\sum_{j=1}^{n_1} (y_{1j} - \bar{y}_1)^2 + \sum_{j=1}^{n_2} (y_{2j} - \bar{y}_2)^2}{n_1 + n_2 - 2}.$$

An attempt is now made to find a linear combination \mathbf{a} of \mathbf{x} , which allows a maximum separation of the sample means \bar{y}_1 and \bar{y}_2 .

Theorem 8.3.4 *The linear combination*

$$\hat{y} = \hat{\mathbf{a}}^\top \mathbf{x} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \mathbf{x} \quad (8.18)$$

maximizes the ratio

$$\frac{(\bar{y}_1 - \bar{y}_2)^2}{s_y^2} = \frac{(\hat{\mathbf{a}}^\top \bar{\mathbf{x}}_1 - \hat{\mathbf{a}}^\top \bar{\mathbf{x}}_2)^2}{\hat{\mathbf{a}}^\top \mathbf{S}_{pooled} \hat{\mathbf{a}}} \quad (8.19)$$

over all possible vectors $\hat{\mathbf{a}}$. \mathbf{S}_{pooled} is defined analogously to the two-group case. The maximum ratio of (8.19) is

$$D^2 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2).$$

Theorem 8.3.5 *An observation \mathbf{x}_0 is assigned to π_1 , if*

$$\hat{y}_0 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \mathbf{x}_0 \geq \hat{m} = \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) \quad (8.20)$$

Otherwise, \mathbf{x}_0 is assigned to π_2 .

This classification rule is shown in Figure 8.4 for $p = 2$. The observations are projected onto a straight line. The direction $\hat{\mathbf{a}}$ of the line is varied until the two groups are maximally separated.

In contrast to rule (8.13), Fisher's classification rule does not explicitly require the assumption of normal distribution of both populations. However, in case of violations from this assumption, we know from the previous results that the rule would not be optimal in terms of minimizing TPM (or ECM). Moreover, we also have to assume that the populations have the same covariance matrix, since a pooled estimate of the covariance is used. It can be seen immediately that Fisher's linear discriminant function in (8.20) is a special case of (8.13). If the prior probabilities and the expected costs for misclassification are the same for rule (8.13), which results from minimizing the expected costs of misclassification, we get exactly (8.20).

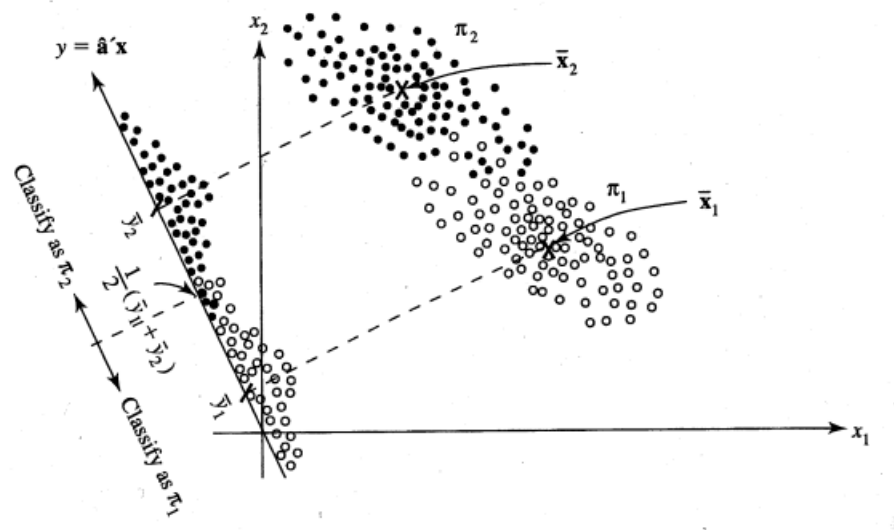


Figure 8.4: Schematic representation of Fisher's classification rule.

- 45 **Extend the ECM to the multi-group case. What is the resulting decision rule if costs are equal? What are the discriminant functions?**

Deviations from the normal distribution of the groups or different covariances could strongly distort the good theoretical properties.

8.4.1 The method to minimize ECM

Let $f_i(\mathbf{x})$ be the density of the observations of population π_i with $i = 1, \dots, g$. Mostly, it is assumed that $f_i(\mathbf{x})$ is the density of a multivariate normal distribution, but this is not a requirement of the following method.

The notation is based on the two-group case. Thus, p_i denotes the prior probability for the population π_i ($i = 1 \dots, g$). Further, let R_k be the space of observations

\mathbf{x} to which the objects from π_k are assigned. $c(k|i)$ is the cost of misclassification when objects from π_i are mistakenly mapped to π_k ($k = 1, \dots, g$), where $k = i$ is of course $c(i|i) = 0$. The probability of this wrong assignment is

$$P(k|i) = P(\mathbf{X} \in R_k | \pi_i) = \int_{R_k} f_i(\mathbf{x}) d\mathbf{x} . \quad (8.21)$$

The conditional expected cost of misclassifying \mathbf{x} from π_1 to π_2, \dots, π_g are analogous to (8.3)

$$ECM(1) = \sum_{k=2}^g P(k|1) c(k|1) . \quad (8.22)$$

These expected costs arise with prior probability p_1 . $ECM(2), \dots, ECM(g)$ are analogously defined, and one obtains by multiplying with the prior probabilities and adding up all contributions, the total expected cost of misclassification (ECM) as

$$ECM = \sum_{i=1}^g p_i ECM(i) = \sum_{i=1}^g p_i \left(\sum_{\substack{k=1 \\ k \neq i}}^g P(k|i) c(k|i) \right) . \quad (8.23)$$

An optimal classification rule should now yield such ranges R_1, \dots, R_g (disjoint and complete decomposition of the sample space Ω) so that (8.23) is minimized.

Theorem 8.4.1 *The areas for minimizing ECM (8.23) are given by assigning \mathbf{x} to the population π_k ($k = 1, \dots, g$) for which the expression*

$$\sum_{\substack{i=1 \\ k \neq i}}^g p_i f_i(\mathbf{x}) c(k|i) \quad (8.24)$$

is minimal.

For simplicity, we assume in the remainder of this chapter that the costs of misclassification are the same for all groups, and thus we can ignore them (or set them equal to 1). Then, minimizing criterion (8.24) corresponds to maximizing the omitted term $p_k f_k(\mathbf{x})$, which simplifies the classification rule to:

An observation \mathbf{x} is assigned to π_k , if

$$p_k f_k(\mathbf{x}) > p_i f_i(\mathbf{x}) \quad \text{for all } i \neq k . \quad (8.25)$$

Note that the above classification rules can only be applied if the prior probabilities, the misclassification costs, and the density functions are known.

8.4.2 Classification in case of the normal distribution

As an important special case, we consider multivariate normally distributed populations π_i with mean $\boldsymbol{\mu}_i$ and covariance matrix $\boldsymbol{\Sigma}_i$,

$$f_i(\mathbf{x}) = (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}_i|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right\} \quad \text{for } i = 1, \dots, g. \quad (8.26)$$

If we now use rule (8.25) or the logarithm of this rule, we get:

An observation \mathbf{x} is assigned to π_k , if

$$\ln p_k f_k(\mathbf{x}) = \ln p_k - \frac{p}{2} \ln(2\pi) - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) = \max_i \ln p_i f_i(\mathbf{x}). \quad (8.27)$$

The constant $(p/2) \ln(2\pi)$ can be omitted in (8.27) since it is the same for all π_i . We therefore define the *quadratic discriminant values* for the i -th population as

$$d_i^Q(\mathbf{x}) = -\frac{1}{2} \ln |\boldsymbol{\Sigma}_i| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \ln p_i \quad \text{for } i = 1, \dots, g. \quad (8.28)$$

Thus one obtains the following classification rule:

An observation \mathbf{x} is assigned to π_k , if:

$$d_k^Q(\mathbf{x}) \quad \text{is the largest of} \quad d_1^Q(\mathbf{x}), \dots, d_g^Q(\mathbf{x}). \quad (8.29)$$

Mostly, the $\boldsymbol{\mu}_i$ and $\boldsymbol{\Sigma}_i$ are unknown and must be estimated. If a “training set” with known class memberships of the observations is available, then the sample means $\bar{\mathbf{x}}_i$ and sample covariance matrices \mathbf{S}_i can be used to estimate the corresponding parameters of the i -th population ($i = 1, \dots, g$). This results in the estimated quadratic discriminant values

$$\hat{d}_i^Q(\mathbf{x}) = -\frac{1}{2} \ln |\mathbf{S}_i| - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_i)^\top \mathbf{S}_i^{-1} (\mathbf{x} - \bar{\mathbf{x}}_i) + \ln p_i, \quad (8.30)$$

with the help of which objects are classified analogous to (8.29).

A simplification results if the covariance matrices of the populations are the same, that is $\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}$ for $i = 1, \dots, g$. In this case, the quadratic discriminant values from (8.28) are

$$d_i^Q(\mathbf{x}) = -\frac{1}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \ln p_i. \quad (8.31)$$

Since the first two terms in (8.31) are the same for all populations, they can be omitted, and we get the *linear discriminant values*

$$d_i(\mathbf{x}) = \boldsymbol{\mu}_i^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \ln p_i \quad \text{for } i = 1, \dots, g. \quad (8.32)$$

One obtains an estimate of the linear discriminant values from the sample by first estimating a pooled covariance matrix

$$\begin{aligned} \mathbf{S}_{pooled} &= \frac{1}{n_1 + \dots + n_g - g} \left((n_1 - 1) \mathbf{S}_1 + \dots + (n_g - 1) \mathbf{S}_g \right) \\ &= \frac{1}{\sum_{i=1}^g n_i - g} \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)(\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)^\top. \end{aligned}$$

This gives the following estimates

$$\hat{d}_i(\mathbf{x}) = \bar{\mathbf{x}}_i^\top \mathbf{S}_{pooled}^{-1} \mathbf{x} - \frac{1}{2} \bar{\mathbf{x}}_i^\top \mathbf{S}_{pooled}^{-1} \bar{\mathbf{x}}_i + \ln p_i . \quad (8.33)$$

The resulting classification rule is:

An observation \mathbf{x} is assigned to π_k , if:

$$\hat{d}_k(\mathbf{x}) \quad \text{is the largest of} \quad \hat{d}_1(\mathbf{x}), \dots, \hat{d}_g(\mathbf{x}) . \quad (8.34)$$

Remark: The expression (8.32) is a linear function of \mathbf{x} . However, one could obtain an analogous rule by ignoring the first term in (8.28), which is the same for all populations in the case of equal covariance matrices. With the corresponding estimated values, one obtains a squared distance

$$D_i^2(\mathbf{x}) = (\mathbf{x} - \bar{\mathbf{x}}_i)^\top \mathbf{S}_{pooled}^{-1} (\mathbf{x} - \bar{\mathbf{x}}_i) , \quad (8.35)$$

and the classification rule is:

$$\text{Assign } \mathbf{x} \text{ to population } \pi_i \text{ for which } -1/2 D_i^2(\mathbf{x}) + \ln p_i \text{ is the largest.} \quad (8.36)$$

This rule is analogous to rule (8.34), it assigns \mathbf{x} to the population that is "closest", with the distance measure penalized with $\ln p_i$. If the a-priori probabilities are not known, one could estimate them by $p_i = 1/g$.

46 What is the idea to extend the two-group Fischer criterion to the multi-group case?

8.4.3 Discriminant analysis by Fisher for the multi-group case

Fisher's discriminant analysis for $g = 2$ can be extended to the multi-group case ($g > 2$), see Rao (1948). For this purpose, we again consider univariate projections of the form $y = \mathbf{a}^T \mathbf{x}$ with $\mathbf{a} \in \mathbb{R}^p$ and $\mathbf{a} \neq \mathbf{0}$, but this time a single projection direction \mathbf{a} will not be enough to describe the solution.

Let $\mathbf{a} \neq \mathbf{0}$ be the projection direction we are looking for. The expected value for population π_i (for $i = 1, \dots, g$) of the random variable $y = \mathbf{a}^T \mathbf{x}$ is then

$$\mu_{i,y} = E(y|\mathbf{x} \in \pi_i) = \mathbf{a}^T E(\mathbf{x}|\mathbf{x} \in \pi_i) = \mathbf{a}^T \boldsymbol{\mu}_i$$

and the variance is

$$\sigma_{i,y}^2 = \text{Var}(y|\mathbf{x} \in \pi_i) = \mathbf{a}^T \text{Cov}(\mathbf{x}|\mathbf{x} \in \pi_i) \mathbf{a} = \mathbf{a}^T \boldsymbol{\Sigma}_i \mathbf{a} .$$

The total weighted average of the populations is denoted by $\bar{\boldsymbol{\mu}} = \sum_{i=1}^g p_i \boldsymbol{\mu}_i$, and the corresponding projection into the univariate space by $\bar{\mu}_y = \mathbf{a}^T \bar{\boldsymbol{\mu}}$.

We now make the assumption that the covariances of all groups are the same, that is $\boldsymbol{\Sigma}_1 = \dots = \boldsymbol{\Sigma}_g = \boldsymbol{\Sigma}$. Then it also applies to the above variance that

$$\sigma_{1,y}^2 = \dots = \sigma_{g,y}^2 = \sigma_y^2 = \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a} .$$

47 How to arrive at the Fischer multi-group solution (variation within and between groups)?

In Fisher's two-group case, an expression $(\bar{y}_1 - \bar{y}_2)^2/s_y^2$ was maximized. In our notation with random variables, this corresponds to $(\mu_{1,y} - \mu_{2,y})^2/\sigma_y^2$. Additionally, we now consider prior probabilities, and the weighted mean is

$$\bar{\mu}_y = p_1\mu_{1,y} + p_2\mu_{2,y} = \mathbf{a}^T(p_1\boldsymbol{\mu}_1 + p_2\boldsymbol{\mu}_2) .$$

Since $p_1 + p_2 = 1$, it is straightforward to see that

$$p_1(\mu_{1,y} - \bar{\mu}_y)^2 + p_2(\mu_{2,y} - \bar{\mu}_y)^2 = p_1p_2(\mu_{1,y} - \mu_{2,y})^2 . \quad (8.37)$$

The latter expression (8.37) should therefore be maximized using the Fisher rule (in terms of variance), and this expression describes the weighted sum of the squared distances of the group means to the total mean.

The generalization of (8.37) to the multi-group case is then immediately obvious; now

$$\frac{\sum_{i=1}^g p_i(\mu_{i,y} - \bar{\mu}_y)^2}{\sigma_y^2} \quad (8.38)$$

should be maximized. The denominator is according to above $\sigma_y^2 = \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}$. If the group covariances are not equal, the resulting classification rule will no longer be optimal. $\boldsymbol{\Sigma}$ would then be best replaced by a pooled version, that is

$$\mathbf{W} = \sum_{i=1}^g p_i \boldsymbol{\Sigma}_i .$$

The matrix \mathbf{W} describes the *variation within groups*.

The numerator of (8.38) can be represented as

$$\sum_{i=1}^g p_i(\mu_{i,y} - \bar{\mu}_y)^2 = \sum_{i=1}^g p_i(\mathbf{a}^T(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}}))^2 = \mathbf{a}^T \mathbf{B} \mathbf{a} ,$$

with the matrix

$$\mathbf{B} = \sum_{i=1}^g p_i(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})^T ,$$

which describes the *variation between groups*.

All in all, the maximization problem (8.38) can be expressed as

$$\frac{\mathbf{a}^T \mathbf{B} \mathbf{a}}{\mathbf{a}^T \mathbf{W} \mathbf{a}} \quad \text{for } \mathbf{a} \in \mathbb{R}^p, \mathbf{a} \neq \mathbf{0} . \quad (8.39)$$

Theorem 8.4.2 *The solution of the maximization problem (8.39) is given by the eigenvectors $\mathbf{a}_1, \dots, \mathbf{a}_l$ of the matrix $\mathbf{W}^{-1}\mathbf{B}$, which must be scaled, so that $\mathbf{a}_j^T \mathbf{W} \mathbf{a}_j = 1$ for $j = 1, \dots, l$. The number l of the strictly positive eigenvalues of the eigenvalue decomposition of $\mathbf{W}^{-1}\mathbf{B}$ is then $l \leq \min(g-1, p)$.*

We can now define the *Fisher discriminant functions* as $y_j = \mathbf{a}_j^T \mathbf{x}$, for $j = 1, \dots, l$, which is the projection of the random variable \mathbf{x} in the direction \mathbf{a}_j . If specific data is available, visualizing the first two discriminant functions is particularly interesting because it represents the projection of the data in which the group means appear best separated. Note that in the case of $g = 3$ groups it holds that $l \leq 2$, regardless of whether p is large or not.

Finally, we also want to get a classification rule. For that, consider the *Fisher discriminant values*

$$d_i^F(\mathbf{x}) = \sum_{j=1}^l (y_j - \mu_{i,y_j})^2 - 2 \log p_i \quad (8.40)$$

for $i = 1, \dots, g$. Here $\mu_{i,y_j} = \mathbf{a}_j^T \boldsymbol{\mu}_i$, and one thus has a measure of the deviation from \mathbf{x} to the i -th group mean in the discriminant space, adjusted with the prior probability (analogous to earlier). Note that here in the discriminant space the distance measure is simply the Euclidean distance. A new observation \mathbf{x} is then assigned to population π_k , if $d_k^F(\mathbf{x})$ is the smallest (!) value of all the values of the groups $d_1^F(\mathbf{x}), \dots, d_g^F(\mathbf{x})$.

By arranging the eigenvectors $\mathbf{a}_1, \dots, \mathbf{a}_l$ in the columns of the matrix \mathbf{A} , we can also write the *Fisher discriminant values* as

$$d_i^F(\mathbf{x}) = \sum_{j=1}^l (\mathbf{a}_j^T (\mathbf{x} - \boldsymbol{\mu}_i))^2 - 2 \log p_i = (\mathbf{x} - \boldsymbol{\mu}_i)^T \mathbf{A} \mathbf{A}^T (\mathbf{x} - \boldsymbol{\mu}_i) - 2 \log p_i ,$$

which corresponds to a (squared) Mahalanobis distance in the original space.