

Notes for ISyE 7401

Alexander Shapiro

Contents

1	Some matrix calculus	3
2	Multivariate normal distribution	5
2.1	Schur complement	8
3	Quadratic forms	9
4	Statistical inference of linear models	10
4.1	Distribution theory	13
4.2	Estimation with linear constraints	15
4.3	Polynomial Regression	17
5	Shrinkage Methods	18
5.1	Ridge Regression	18
5.2	Lasso method	20
6	Elements of large samples theory	21
7	Exponential family of distributions	24
8	Point estimation	25
8.1	Maximum likelihood method	25
8.1.1	Asymptotic distribution of the ML estimators	26
8.2	Cramér - Rao lower bound	30
8.3	Best unbiased estimators	32
9	Hypotheses testing	35
9.1	Likelihood Ratio Test	37
9.2	Testing equality constraints	38
10	Multinomial distribution	39
11	Logistic regression	41
12	Generalized linear models	43
13	Classification problem	45
13.1	Classification with normally distributed populations	46
13.1.1	An optimization problem	47
13.2	Fisher discriminant analysis	48
13.3	Several populations	48
13.3.1	Mahalanobis distance	49
13.4	Bayes and Logistic Regression classifiers	49
14	Support Vector Machines	50

15 Principal Components Analysis	53
15.1 Derivatives of eigenvalues and eigenvectors	54
15.2 Elements of matrix calculus	56
15.3 Asymptotics of PCA	57
15.4 Singular value decomposition	58
16 Factor analysis model	59
17 Kernel PCA	61
18 Correlation analysis	62
18.1 Partial correlation	62
18.2 Canonical correlation analysis	62
19 Gaussian Mixture Models	64
20 Von Mises statistical functionals	64
21 Bootstrap	68
21.1 Jackknife bias estimation.	68
21.2 Bootstrap method	68
22 Robust statistics	69
22.1 Quantile regression	69
23 Bayes estimators	70
23.1 Bayesian decisions	71
24 Spherical and elliptical distributions	75
24.1 Multivariate cumulants	77
25 Wishart distribution	79
25.1 Hotelling's T^2 statistic	82
26 Spatial statistics	83

1 Some matrix calculus

For an $m \times n$ matrix \mathbf{A} we denote by \mathbf{A}' its transpose $n \times m$ matrix. Unless stated otherwise vectors $\mathbf{a} = (a_1, \dots, a_m)'$ are assumed to be column vectors. If \mathbf{A} and \mathbf{B} are two matrices such that their product \mathbf{AB} is well defined, then the transpose of \mathbf{AB} is $\mathbf{B}'\mathbf{A}'$, i.e., $(\mathbf{AB})' = \mathbf{B}'\mathbf{A}'$. Trace of a square $m \times m$ matrix \mathbf{A} is defined as the sum of its diagonal elements, i.e., $\text{tr}(\mathbf{A}) = a_{11} + \dots + a_{mm}$. It has the following important property. Let \mathbf{A} and \mathbf{B} be two matrices such that their product \mathbf{AB} is well defined. Then

$$\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA}). \quad (1.1)$$

In particular, if $\mathbf{a} = (a_1, \dots, a_m)'$ is an $m \times 1$ vector, then \mathbf{aa}' is an $m \times m$ matrix, its trace $\text{tr}(\mathbf{aa}') = \sum_{i=1}^m a_i^2 = \mathbf{a}'\mathbf{a}$.

Let \mathbf{A} be an $m \times m$ matrix. We denote by $|\mathbf{A}|$ the determinant of \mathbf{A} . Matrix \mathbf{A} is nonsingular (invertible) if and only if (iff) $|\mathbf{A}| \neq 0$. It is said that λ is an eigenvalue of \mathbf{A} if there is an $m \times 1$ vector $\mathbf{e} \neq \mathbf{0}$ such that $\mathbf{A}\mathbf{e} = \lambda\mathbf{e}$. It follows that $(\mathbf{A} - \lambda\mathbf{I}_m)\mathbf{e} = \mathbf{0}$, where \mathbf{I}_m is the $m \times m$ identity matrix. Thus matrix $(\mathbf{A} - \lambda\mathbf{I}_m)$ is singular, and hence its determinant $|\mathbf{A} - \lambda\mathbf{I}_m| = 0$. Consider¹ $p(\lambda) := |\mathbf{A} - \lambda\mathbf{I}_m|$. This is a polynomial of degree m and hence has m roots which are eigenvalues of matrix \mathbf{A} . Therefore matrix \mathbf{A} has m eigenvalues some of which can be complex numbers. Suppose now that matrix \mathbf{A} is symmetric, i.e., $\mathbf{A}' = \mathbf{A}$. Then it has m real valued eigenvalues $\lambda_1 \geq \dots \geq \lambda_m$ and a corresponding set of eigenvectors $\mathbf{e}_1, \dots, \mathbf{e}_m$ such that

$$\mathbf{A}\mathbf{e}_i = \lambda_i\mathbf{e}_i, \quad i = 1, \dots, m. \quad (1.2)$$

The eigenvectors can be chosen in such a way that $\mathbf{e}'_i\mathbf{e}_j = 0$ for $i \neq j$ and $\mathbf{e}'_i\mathbf{e}_i = 1$ for $i = 1, \dots, m$, i.e., these eigenvectors are orthogonal to each other and of length one. In that case we say the eigenvectors are *orthonormal*.

Consider the $m \times m$ matrix $\mathbf{T} = [\mathbf{e}_1, \dots, \mathbf{e}_m]$ whose columns are formed from a set of orthonormal eigenvectors. Matrix \mathbf{T} has the following property $\mathbf{T}'\mathbf{T} = \mathbf{I}_m$ and $\mathbf{T}\mathbf{T}' = \mathbf{I}_m$. Such matrices are called *orthogonal*. Equations (1.2) can be written in the form $\mathbf{A}\mathbf{T} = \mathbf{T}\mathbf{\Lambda}$, where $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_m)$ is the diagonal matrix. By multiplying both sides of this matrix equation by \mathbf{T}' we obtain

$$\mathbf{A} = \mathbf{T}\mathbf{\Lambda}\mathbf{T}' = \sum_{i=1}^m \lambda_i \mathbf{e}_i \mathbf{e}'_i. \quad (1.3)$$

The representation (1.3) is called *spectral decomposition* of matrix \mathbf{A} . It also follows that $\mathbf{T}'\mathbf{A}\mathbf{T} = \mathbf{\Lambda}$, and that $\text{tr}(\mathbf{A}) = \text{tr}(\mathbf{\Lambda}) = \lambda_1 + \dots + \lambda_m$, and $\mathbf{A}^{-1} = \mathbf{T}\mathbf{\Lambda}^{-1}\mathbf{T}'$, provided that all $\lambda_i \neq 0$, $i = 1, \dots, m$.

It is said that matrix \mathbf{A} is *positive semidefinite* if $\mathbf{x}'\mathbf{A}\mathbf{x} \geq 0$ for any $\mathbf{x} \in \mathbb{R}^m$, and it is said that \mathbf{A} is *positive definite* if $\mathbf{x}'\mathbf{A}\mathbf{x} > 0$ for any $\mathbf{x} \neq \mathbf{0}$. By using (1.3) we can write

$$\mathbf{x}'\mathbf{A}\mathbf{x} = \mathbf{x}'\mathbf{T}\mathbf{\Lambda}\mathbf{T}'\mathbf{x} = \mathbf{y}'\mathbf{\Lambda}\mathbf{y} = \sum_{i=1}^m \lambda_i y_i^2,$$

where $\mathbf{y} = \mathbf{T}'\mathbf{x}$. Note that $\mathbf{y}'\mathbf{y} = \mathbf{x}'\mathbf{T}\mathbf{T}'\mathbf{x} = \mathbf{x}'\mathbf{x}$. It follows that matrix \mathbf{A} is positive semidefinite iff all its eigenvalues are nonnegative, and is positive definite iff all its eigenvalues are positive.

We can define a function of symmetric matrix \mathbf{A} by considering a function of its eigenvalues. For example if matrix \mathbf{A} is positive semidefinite and hence all its eigenvalues are

¹Sometimes we write ':= ' meaning 'equal by definition'.

nonnegative we can define $\mathbf{A}^{1/2} = \mathbf{T}\mathbf{\Lambda}^{1/2}\mathbf{T}'$, where $\mathbf{\Lambda}^{1/2} = \text{diag}(\lambda_1^{1/2}, \dots, \lambda_m^{1/2})$. The so defined matrix $\mathbf{A}^{1/2}$ is symmetric positive semidefinite and $(\mathbf{A}^{1/2})^2 = \mathbf{T}\mathbf{\Lambda}^{1/2}\mathbf{T}'\mathbf{T}\mathbf{\Lambda}^{1/2}\mathbf{T}' = \mathbf{A}$, since $\mathbf{T}'\mathbf{T} = \mathbf{I}_m$. Similarly if \mathbf{A} is positive definite and hence all its eigenvalues are positive, we can define $\mathbf{A}^{-1/2} = \mathbf{T}\mathbf{\Lambda}^{-1/2}\mathbf{T}'$. Matrix $\mathbf{A}^{-1/2}$ is symmetric positive definite and $(\mathbf{A}^{-1/2})^2 = \mathbf{T}\mathbf{\Lambda}^{-1/2}\mathbf{T}'\mathbf{T}\mathbf{\Lambda}^{-1/2}\mathbf{T}' = \mathbf{A}^{-1}$.

Let \mathbf{A} and \mathbf{B} be two $m \times m$ symmetric matrices. Then $(\mathbf{AB})' = \mathbf{BA}$, so the product matrix is not symmetric unless $\mathbf{AB} = \mathbf{BA}$. Suppose that \mathbf{A} is positive semidefinite, then matrix $\mathbf{A}^{1/2}\mathbf{BA}^{1/2}$ is symmetric. Let \mathbf{e} be an eigenvector and λ the corresponding eigenvalue of $\mathbf{A}^{1/2}\mathbf{BA}^{1/2}$, i.e., $\mathbf{A}^{1/2}\mathbf{BA}^{1/2}\mathbf{e} = \lambda\mathbf{e}$. Multiplying both sides of this equation by $\mathbf{A}^{1/2}$ we obtain $\mathbf{ABA}^{1/2}\mathbf{e} = \lambda\mathbf{A}^{1/2}\mathbf{e}$. That is, $\mathbf{A}^{1/2}\mathbf{e}$ is the corresponding eigenvector and λ is the eigenvalue of matrix \mathbf{AB} . This shows that although \mathbf{AB} is not symmetric, it has real valued eigenvectors and eigenvalues. Moreover, if \mathbf{B} is positive semidefinite, then $\mathbf{A}^{1/2}\mathbf{BA}^{1/2}$ is positive semidefinite and hence all eigenvalues of \mathbf{AB} are nonnegative, and if both \mathbf{A} and \mathbf{B} are positive definite matrices, then $\mathbf{A}^{1/2}\mathbf{BA}^{1/2}$ is positive definite and hence all eigenvalues of \mathbf{AB} are positive.

Random vectors. Consider an $m \times 1$ random vector $\mathbf{X} = (X_1, \dots, X_m)'$. Its expected value $\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}]$ is defined as $\mathbb{E}[\mathbf{X}] = (\mathbb{E}[X_1], \dots, \mathbb{E}[X_m])'$, i.e., the expectation is taken componentwise. Similarly expectation of a random matrix is taken componentwise. Sometimes we write $\boldsymbol{\mu}_X$ to emphasize that this is mean vector of \mathbf{X} . The $m \times m$ covariance matrix of \mathbf{X} is

$$\boldsymbol{\Sigma} := \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})'] = \mathbb{E}[\mathbf{X}\mathbf{X}'] - \boldsymbol{\mu}\boldsymbol{\mu}'.$$

The (i, j) -component of $\boldsymbol{\Sigma}$ is the covariance $\text{Cov}(X_i, X_j)$, $i, j = 1, \dots, m$.

Covariance matrix $\boldsymbol{\Sigma}$ has the following properties. It is symmetric, i.e., $\boldsymbol{\Sigma}' = \boldsymbol{\Sigma}$. Consider a (deterministic) $k \times m$ matrix \mathbf{A} and $k \times 1$ random vector $\mathbf{Y} = \mathbf{A}\mathbf{X}$. Then

$$\boldsymbol{\mu}_Y = \mathbb{E}[\mathbf{Y}] = \mathbb{E}[\mathbf{A}\mathbf{X}] = \mathbf{A}\mathbb{E}[\mathbf{X}] = \mathbf{A}\boldsymbol{\mu}_X.$$

In particular, if $k = 1$ and $Y = \mathbf{a}'\mathbf{X} = a_1X_1 + \dots + a_mX_m$, where $\mathbf{a} = (a_1, \dots, a_m)'$, then $\mathbb{E}[Y] = \mathbf{a}'\boldsymbol{\mu}_X$. Now

$$\boldsymbol{\Sigma}_Y = \mathbb{E}[\mathbf{Y}\mathbf{Y}'] - \boldsymbol{\mu}_Y\boldsymbol{\mu}_Y' = \mathbb{E}[\mathbf{A}\mathbf{X}\mathbf{X}'\mathbf{A}'] - \mathbf{A}\boldsymbol{\mu}_X\boldsymbol{\mu}_X'\mathbf{A}'.$$

Since $\mathbb{E}[\mathbf{A}\mathbf{X}\mathbf{X}'\mathbf{A}'] = \mathbf{A}\mathbb{E}[\mathbf{X}\mathbf{X}']\mathbf{A}'$ it follows that

$$\boldsymbol{\Sigma}_Y = \mathbf{A}\boldsymbol{\Sigma}_X\mathbf{A}'. \quad (1.4)$$

In particular, if $\mathbf{A} = \mathbf{a}' = (a_1, \dots, a_m)'$ is a row vector, then

$$\text{Var}(\mathbf{a}'\mathbf{X}) = \mathbf{a}'\boldsymbol{\Sigma}\mathbf{a} = \sum_{i,j=1}^m \sigma_{ij}a_ia_j, \quad (1.5)$$

where $\sigma_{ij} = \text{Cov}(X_i, X_j)$ is the (i, j) -component of covariance matrix $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_X$. Since variance of a random variable is always nonnegative, it follows that $\mathbf{a}'\boldsymbol{\Sigma}\mathbf{a} \geq 0$ for any $m \times 1$ vector \mathbf{a} . Therefore covariance matrix $\boldsymbol{\Sigma}$ is positive semidefinite. If moreover $\boldsymbol{\Sigma}$ is nonsingular (invertible), then it is positive definite.

Recall that matrix $\boldsymbol{\Sigma}$ is positive definite iff $\mathbf{a}'\boldsymbol{\Sigma}\mathbf{a} > 0$ for all $\mathbf{a} \neq \mathbf{0}$. If $\mathbf{a}'\boldsymbol{\Sigma}\mathbf{a} = 0$ for some $\mathbf{a} \neq \mathbf{0}$, then this means that $\text{Var}(\mathbf{a}'\mathbf{X}) = 0$ and hence $Y = \mathbf{a}'\mathbf{X}$ is constant. In turn this means that random variables $X_1 - \mu_1, \dots, X_m - \mu_m$ are linearly dependent. Therefore $\boldsymbol{\Sigma}$ is positive

definite iff variables $X_1 - \mu_1, \dots, X_m - \mu_m$ are linearly independent.

As an example let us compute expectation of $\mathbf{X}'\mathbf{A}\mathbf{X} = \sum_{i,j=1}^m a_{ij}X_iX_j$, where \mathbf{A} is an $m \times m$ matrix. Note that using property (1.1) we can write $\mathbf{X}'\mathbf{A}\mathbf{X} = \text{tr}(\mathbf{X}'\mathbf{A}\mathbf{X}) = \text{tr}(\mathbf{A}\mathbf{X}\mathbf{X}')$. Also $\mathbb{E}[\text{tr}(\mathbf{A}\mathbf{X}\mathbf{X}')] = \text{tr}(\mathbb{E}[\mathbf{A}\mathbf{X}\mathbf{X}'])$ and hence

$$\mathbb{E}[\mathbf{X}'\mathbf{A}\mathbf{X}] = \text{tr}(\mathbf{A}\mathbb{E}[\mathbf{X}\mathbf{X}']) = \text{tr}(\mathbf{A}(\boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}')) = \text{tr}(\mathbf{A}\boldsymbol{\Sigma}) + \text{tr}(\mathbf{A}\boldsymbol{\mu}\boldsymbol{\mu}') = \text{tr}(\mathbf{A}\boldsymbol{\Sigma}) + \boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu}. \quad (1.6)$$

2 Multivariate normal distribution

Recall that a random variable X has normal distribution with mean μ and variance σ^2 , denoted $X \sim \mathcal{N}(\mu, \sigma^2)$, if its probability density function (pdf) is

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

Now let X_1, \dots, X_m be an iid sequence² of standard normal variables, i.e., $X_i \sim \mathcal{N}(0, 1)$, $i = 1, \dots, m$, and these random variables are independent of each other. Then the pdf of random vector $\mathbf{X} = (X_1, \dots, X_m)'$ is

$$f_{\mathbf{X}}(\mathbf{x}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} e^{-\frac{x_i^2}{2}} = \frac{1}{(2\pi)^{m/2}} e^{-\frac{x_1^2 + \dots + x_m^2}{2}} = \frac{1}{(2\pi)^{m/2}} \exp(-\mathbf{x}'\mathbf{x}/2).$$

Consider $\mathbf{Y} = \mathbf{A}\mathbf{X}$, where \mathbf{A} is an $m \times m$ nonsingular matrix. Note that $\mathbf{X} = \mathbf{A}^{-1}\mathbf{Y}$. Then the pdf of \mathbf{Y} is

$$f_{\mathbf{Y}}(\mathbf{y}) = f_{\mathbf{X}}(\mathbf{A}^{-1}\mathbf{y}) |\mathbf{A}^{-1}| = \frac{1}{(2\pi)^{m/2} |\mathbf{A}|} \exp(-\mathbf{y}'\mathbf{A}'^{-1}\mathbf{A}^{-1}\mathbf{y}/2) = \frac{1}{(2\pi)^{m/2} |\boldsymbol{\Sigma}_{\mathbf{Y}}|^{1/2}} \exp(-\mathbf{y}'\boldsymbol{\Sigma}_{\mathbf{Y}}^{-1}\mathbf{y}/2).$$

Recall that $|\mathbf{A}|$ denotes determinant of (square) matrix \mathbf{A} . We used the following properties in the above derivations: $|\mathbf{A}^{-1}| = |\mathbf{A}|^{-1}$, $\mathbf{A}'^{-1}\mathbf{A}^{-1} = (\mathbf{A}\mathbf{A}')^{-1}$, and $\boldsymbol{\Sigma}_{\mathbf{Y}} = \mathbf{A}\boldsymbol{\Sigma}_{\mathbf{X}}\mathbf{A}' = \mathbf{A}\mathbf{A}'$ since $\boldsymbol{\Sigma}_{\mathbf{X}} = \mathbf{I}_m$ is the identity matrix, $|\boldsymbol{\Sigma}_{\mathbf{Y}}| = |\mathbf{A}\mathbf{A}'| = |\mathbf{A}||\mathbf{A}'| = |\mathbf{A}|^2$. Note also that $\mathbb{E}[\mathbf{Y}] = \mathbf{A}\boldsymbol{\mu}_{\mathbf{X}} = \mathbf{0}$.

Finally consider $\mathbf{Y} = \mathbf{A}\mathbf{X} + \boldsymbol{\mu}$. The pdf of this random vector is

$$f_{\mathbf{Y}}(\mathbf{y}) = \frac{1}{(2\pi)^{m/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{(\mathbf{y} - \boldsymbol{\mu})'\boldsymbol{\Sigma}^{-1}(\mathbf{y} - \boldsymbol{\mu})}{2}\right\}. \quad (2.1)$$

If random vector \mathbf{Y} has pdf of the form (2.1), where $\boldsymbol{\Sigma}$ is a symmetric positive definite matrix, then it is said that \mathbf{Y} has multivariate normal distribution, denoted $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Sometimes we write this as $\mathcal{N}_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ to emphasize dimension m of random vector \mathbf{Y} . Note that $\boldsymbol{\mu}$ is the $m \times 1$ mean vector and $\boldsymbol{\Sigma}$ is the $m \times m$ covariance matrix of \mathbf{Y} .

Suppose that $\mathbf{X} \sim \mathcal{N}_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is partitioned $\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix}$, where \mathbf{X}_1 and \mathbf{X}_2 are subvectors of \mathbf{X} of the respective dimensions $m_1 \times 1$ and $m_2 \times 1$, with $m_1 + m_2 = m$. The corresponding partitioning of $\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}$ and $\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$. Note that $\boldsymbol{\Sigma}_{21} = \boldsymbol{\Sigma}'_{12}$ since $\boldsymbol{\Sigma}$ is symmetric.

²A sequence X_1, \dots, X_m of random variables is said to be iid (independent identically distributed), if these random variables are independent of each other and have the same probability distribution.

Suppose further that $\Sigma_{12} = \mathbf{0}$ and hence $\Sigma_{21} = \Sigma'_{12} = \mathbf{0}$, i.e., matrix Σ is block diagonal. Then $|\Sigma| = |\Sigma_{11}||\Sigma_{22}|$ and $\Sigma^{-1} = \begin{bmatrix} \Sigma_{11}^{-1} & \mathbf{0} \\ \mathbf{0} & \Sigma_{22}^{-1} \end{bmatrix}$, and hence

$$f_X(\mathbf{x}) = f_{X_1}(\mathbf{x}_1)f_{X_2}(\mathbf{x}_2),$$

where $f_X(\cdot)$ is the pdf of \mathbf{X} and $f_{X_1}(\cdot)$ and $f_{X_2}(\cdot)$ are pdfs of \mathbf{X}_1 and \mathbf{X}_2 , respectively. It follows that random vectors \mathbf{X}_1 and \mathbf{X}_2 are independent. That is, for multivariate normal distribution “independent” and “uncorrelated” are equivalent.

Moment generating function of a random variable X is defined as $M(t) := \mathbb{E}[e^{tX}]$. Since $e^0 = 1$, it follows that $M(0) = 1$. Note that it can happen that $M(t) = +\infty$ for any $t \neq 0$. Two random variables X and Y have the same distribution if their moment generating functions $M_X(t)$ and $M_Y(t)$ are equal to each other for all t in some neighborhood of zero, provided these moment generating functions are finite valued in that neighborhood.

Similarly moment generating function of a random vector $\mathbf{X} = (X_1, \dots, X_m)'$ is defined as

$$M(\mathbf{t}) := \mathbb{E}[e^{t_1 X_1 + \dots + t_m X_m}] = \mathbb{E}[\exp(\mathbf{t}'\mathbf{X})].$$

If $M_X(\mathbf{t})$ is finite valued in a neighborhood of $\mathbf{0} \in \mathbb{R}^m$, then it is differentiable in that neighborhood. Consider $m \times 1$ vector $\partial M_X(\mathbf{t})/\partial \mathbf{t} = (\partial M_X(\mathbf{t})/\partial t_1, \dots, \partial M_X(\mathbf{t})/\partial t_m)'$ of first order partial derivatives, and $m \times m$ matrix of second order³ partial derivatives $\partial^2 M_X(\mathbf{t})/\partial \mathbf{t} \partial \mathbf{t}'$ with (i, j) -element $\partial^2 M_X(\mathbf{t})/\partial t_i \partial t_j$, $i, j = 1, \dots, m$. Then the expectation and differentiation operations can be interchanged (see Remark 8.1) and

$$\partial M_X(\mathbf{t})/\partial \mathbf{t} \Big|_{\mathbf{t}=\mathbf{0}} = \mathbb{E}[\partial \exp(\mathbf{t}'\mathbf{X})/\partial \mathbf{t} \Big|_{\mathbf{t}=\mathbf{0}}] = \mathbb{E}[\mathbf{X} \exp(\mathbf{t}'\mathbf{X}) \Big|_{\mathbf{t}=\mathbf{0}}] = \mathbb{E}[\mathbf{X}].$$

Similarly the Hessian matrix

$$\partial^2 M_X(\mathbf{t})/\partial \mathbf{t} \partial \mathbf{t}' \Big|_{\mathbf{t}=\mathbf{0}} = \mathbb{E}[\mathbf{X} \mathbf{X}'].$$

Let us compute the moment generating function of $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{X})$. For standard normal random variable $X \sim \mathcal{N}(0, 1)$ we have

$$M(t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{tx} e^{-x^2/2} dx = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-(x-t)^2/2} e^{t^2/2} dx = e^{t^2/2} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} e^{-x^2/2} dx = e^{t^2/2}.$$

Let $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_m)$ and hence components $X_i \sim \mathcal{N}(0, 1)$ of \mathbf{X} are independent. Thus

$$M(\mathbf{t}) = \mathbb{E}[e^{t_1 X_1 + \dots + t_m X_m}] = \mathbb{E}[e^{t_1 X_1} \times \dots \times e^{t_m X_m}] = \prod_{i=1}^m \mathbb{E}[e^{t_i X_i}] = \prod_{i=1}^m e^{t_i^2/2} = \exp(\mathbf{t}'\mathbf{t}/2).$$

Consider now $\mathbf{Y} = \mathbf{A}\mathbf{X} + \boldsymbol{\mu}$. Since $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_m)$ we have that $\mathbb{E}[\mathbf{Y}] = \boldsymbol{\mu}$ and the covariance matrix of \mathbf{Y} is $\Sigma = \mathbf{A}\mathbf{A}'$. Then

$$\begin{aligned} M_Y(\mathbf{t}) &= \mathbb{E}[\exp(\mathbf{t}'(\mathbf{A}\mathbf{X} + \boldsymbol{\mu}))] = \mathbb{E}[\exp(\mathbf{t}'\boldsymbol{\mu}) \exp(\mathbf{t}'\mathbf{A}\mathbf{X})] = \exp(\mathbf{t}'\boldsymbol{\mu}) \mathbb{E}[\exp((\mathbf{A}'\mathbf{t})'\mathbf{X})] \\ &= \exp(\mathbf{t}'\boldsymbol{\mu}) M_X(\mathbf{A}'\mathbf{t}) = \exp(\mathbf{t}'\boldsymbol{\mu}) \exp(\mathbf{t}'\mathbf{A}\mathbf{A}'\mathbf{t}/2) = \exp(\mathbf{t}'\boldsymbol{\mu} + \mathbf{t}'\Sigma\mathbf{t}/2). \end{aligned}$$

That is, for $\mathbf{Y} \sim \mathcal{N}_m(\boldsymbol{\mu}, \Sigma)$ its moment generating function is finite valued for any $m \times 1$ vector \mathbf{t} , and

$$M_Y(\mathbf{t}) = \exp(\mathbf{t}'\boldsymbol{\mu} + \mathbf{t}'\Sigma\mathbf{t}/2). \quad (2.2)$$

³Matrix of second order partial derivatives is called Hessian matrix.

Now let $\mathbf{X} \sim \mathcal{N}_m(\boldsymbol{\mu}, \boldsymbol{\Sigma}_X)$ and $\mathbf{Y} = \mathbf{A}\mathbf{X} + \boldsymbol{\eta}$, where \mathbf{A} is a $k \times m$ matrix and $\boldsymbol{\eta}$ is $k \times 1$ vector. Then

$$M_Y(\mathbf{t}) = \exp(\mathbf{t}'(\mathbf{A}\mathbf{X} + \boldsymbol{\eta})) = \exp(\mathbf{t}'\boldsymbol{\eta})M_X(\mathbf{A}'\mathbf{t}) = \exp(\mathbf{t}'(\mathbf{A}\boldsymbol{\mu} + \boldsymbol{\eta})) \exp(\mathbf{t}'\mathbf{A}\boldsymbol{\Sigma}_X\mathbf{A}'\mathbf{t}/2).$$

That is, the moment generating function of \mathbf{Y} is the same as the moment generating function of multivariate normal with mean $\mathbf{A}\boldsymbol{\mu} + \boldsymbol{\eta}$ and covariance matrix $\mathbf{A}\boldsymbol{\Sigma}_X\mathbf{A}'$. It follows that \mathbf{Y} has multivariate normal distribution with mean $\boldsymbol{\mu}_Y = \mathbf{A}\boldsymbol{\mu} + \boldsymbol{\eta}$ and covariance matrix $\boldsymbol{\Sigma}_Y = \mathbf{A}\boldsymbol{\Sigma}_X\mathbf{A}'$. In particular, marginal distribution of every subvector of \mathbf{X} is multivariate normal.

A delicate point of the above result is that the covariance matrix $\mathbf{A}\boldsymbol{\Sigma}_X\mathbf{A}'$ of \mathbf{Y} should be non-singular, i.e. *positive* definite, in order for its density function $f_Y(\mathbf{y})$ to be well defined. Since the covariance matrix $\boldsymbol{\Sigma}$ of \mathbf{X} is positive definite, the matrix $\mathbf{A}\boldsymbol{\Sigma}_X\mathbf{A}'$ is nonsingular iff the $k \times m$ matrix \mathbf{A} has rank k . For example, if $k > m$, then $\text{rank}(\mathbf{A}) \leq m < k$ and hence $\mathbf{A}\boldsymbol{\Sigma}_X\mathbf{A}'$ is singular.

It follows that random vector \mathbf{X} has multivariate normal distribution iff $Y = \mathbf{a}'\mathbf{X}$ is normally distributed for any vector $\mathbf{a} \neq \mathbf{0}$. Indeed, if \mathbf{X} has normal distribution, then as it was shown above $\mathbf{a}'\mathbf{X}$ is normally distributed. Conversely, suppose that $\mathbf{a}'\mathbf{X}$ is normally distributed for any $\mathbf{a} \neq \mathbf{0}$. Consider $Y := \mathbf{t}'\mathbf{X}$ for $\mathbf{t} \neq \mathbf{0}$. We have that $\mu_Y = \mathbf{t}'\boldsymbol{\mu}_X$ and $\sigma_Y^2 = \mathbf{t}'\boldsymbol{\Sigma}_X\mathbf{t}$. Moreover since Y has normal distribution its moment generating function $M_Y(t) = \exp(\mu_Y t + \sigma_Y^2 t^2/2)$. It follows that

$$M_X(\mathbf{t}) = \mathbb{E}[\exp(\mathbf{t}'\mathbf{X})] = M_Y(1) = \exp(\mu_Y + \sigma_Y^2/2) = \exp(\mathbf{t}'\boldsymbol{\mu}_X + \mathbf{t}'\boldsymbol{\Sigma}_X\mathbf{t}/2).$$

That is, the moment generating function of \mathbf{X} has the form of normal distribution (see equation (2.2)). It follows that \mathbf{X} has normal distribution. \square

Conditional normal distribution. Suppose that $\mathbf{X} \sim \mathcal{N}_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is partitioned $\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix}$ with the corresponding partitioning of $\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}$ and $\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$. We want to compute the conditional distribution of \mathbf{X}_1 given $\mathbf{X}_2 = \mathbf{x}_2$. Consider

$$\mathbf{Y} := \mathbf{X}_1 - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\mathbf{X}_2 = [\mathbf{I}_{m_1}, -\boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}] \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix}.$$

Note that vector $(\mathbf{Y}', \mathbf{X}_2')$ has multivariate normal distribution. Moreover $\mathbf{X}_2 = [\mathbf{0}, \mathbf{I}_{m_2}] \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix}$ and

$$\text{Cov}[\mathbf{Y}, \mathbf{X}_2] = [\mathbf{I}_{m_1}, -\boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}] \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ \mathbf{I}_{m_2} \end{bmatrix} = [\mathbf{I}_{m_1}, -\boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}] \begin{bmatrix} \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{22} \end{bmatrix} = \mathbf{0}.$$

It follows that \mathbf{Y} and \mathbf{X}_2 are uncorrelated and hence independent. Since $\mathbf{X}_1 = \mathbf{Y} + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\mathbf{X}_2$ it follows that the conditional distribution of \mathbf{X}_1 given $\mathbf{X}_2 = \mathbf{x}_2$ is the same as the distribution of $\mathbf{Y} + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\mathbf{x}_2$. Now \mathbf{Y} has multivariate normal distribution with mean

$$\mathbb{E}[\mathbf{Y}] = \boldsymbol{\mu}_1 - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\mu}_2$$

and covariance matrix

$$\boldsymbol{\Sigma}_Y = [\mathbf{I}_{m_1}, -\boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}] \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{I}_{m_1} \\ -\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21} \end{bmatrix} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}.$$

That is the conditional distribution of \mathbf{X}_1 given $\mathbf{X}_2 = \mathbf{x}_2$ is multivariate normal

$$\mathcal{N}_{m_1}(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}). \quad (2.3)$$

Note that the conditional covariance matrix $\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$ is given by the Schur complement of $\boldsymbol{\Sigma}$ (see equation (2.5) below).

Chi-square, t and F distributions. There are three important distributions derived from the normal distribution. Note that if $X \sim \mathcal{N}(0, 1)$, then $\mathbb{E}[X^2] = \text{Var}(X) = 1$ and, using integration by parts,

$$\mathbb{E}[X^4] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} x^4 e^{-x^2/2} dx = \frac{3}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} x^2 e^{-x^2/2} dx = 3\text{Var}(X) = 3. \quad (2.4)$$

Hence $\text{Var}(X^2) = \mathbb{E}[X^4] - 1 = 2$.

Let Z_1, \dots, Z_m be an iid sequence of standard normal random variables. Then $Y := Z_1^2 + \dots + Z_m^2$ has chi-square distribution with m degrees of freedom, denoted $Y \sim \chi_m^2$. The expected value of Y is $\mathbb{E}[Y] = \mathbb{E}[Z_1^2] + \dots + \mathbb{E}[Z_m^2] = m$ and variance $\text{Var}(Y) = \text{Var}(Z_1^2) + \dots + \text{Var}(Z_m^2) = 2m$. By the Law of Large Numbers, Y/m tends in probability to 1, and by the Central Limit Theorem, $m^{-1/2}(Y - m)$ tends in distribution to normal $N(0, 2)$, as $m \rightarrow \infty$.

The t distribution with m degrees of freedom is defined as distribution of $T = \frac{Z}{\sqrt{W/m}}$, where $Z \sim \mathcal{N}(0, 1)$ and $W \sim \chi_m^2$ are independent random variables, denoted $T \sim t_m$. Since for large m , W/m becomes close to one, critical values of t -statistic are close to the respective standard normal critical values when the degrees of freedom are large.

The F distribution with k and m degrees of freedom is defined as distribution of $F = \frac{V/k}{W/m}$, where $V \sim \chi_k^2$ and $W \sim \chi_m^2$ are independent random variables, denoted $F \sim F_{k,m}$. It follows from the above definitions that if $T \sim t_m$, then $T^2 \sim F_{1,m}$.

2.1 Schur complement

Consider $(n + m) \times (n + m)$ matrix

$$M = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix},$$

where $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$ are matrices of respective dimensions $n \times n$, $n \times m$, $m \times n$, $m \times m$. Suppose that \mathbf{D} is invertible (nonsingular). Then

$$\begin{bmatrix} \mathbf{I}_n & -\mathbf{B}\mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I}_m \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ -\mathbf{D}^{-1}\mathbf{C} & \mathbf{I}_m \end{bmatrix} = \begin{bmatrix} \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}. \quad (2.5)$$

The matrix $\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}$ is called the Schur complement of M with respect to \mathbf{D} .

Note that

$$\begin{bmatrix} \mathbf{I}_n & -\mathbf{B}\mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I}_m \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{I}_n & \mathbf{B}\mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I}_m \end{bmatrix}$$

and

$$\begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ -\mathbf{D}^{-1}\mathbf{C} & \mathbf{I}_m \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{D}^{-1}\mathbf{C} & \mathbf{I}_m \end{bmatrix},$$

and determinants of these matrices equal one. Hence it follows from (2.5) that

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I}_n & \mathbf{B}\mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I}_m \end{bmatrix} \begin{bmatrix} \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{D}^{-1}\mathbf{C} & \mathbf{I}_m \end{bmatrix}.$$

This implies the following formula for the determinant of matrix \mathbf{M} ,

$$|\mathbf{M}| = |\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}| |\mathbf{D}|. \quad (2.6)$$

Also it follows that matrix \mathbf{M} is invertible iff the matrix $\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}$ is invertible (recall that it is assumed that \mathbf{D} is invertible), in which case

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{I}_n & \mathbf{0} \\ -\mathbf{D}^{-1}\mathbf{C} & \mathbf{I}_m \end{bmatrix} \begin{bmatrix} (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{I}_n & -\mathbf{B}\mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I}_m \end{bmatrix}.$$

Using the above equation it is possible to compute

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & -(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} \end{bmatrix}. \quad (2.7)$$

3 Quadratic forms

In this section we discuss distribution of quadratic forms $Q = \mathbf{X}'\mathbf{A}\mathbf{X}$, where $\mathbf{X} \sim \mathcal{N}_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and \mathbf{A} is an $m \times m$ symmetric (deterministic) matrix. Recall that the expected value of $\mathbf{X}'\mathbf{A}\mathbf{X}$ was computed in equation (1.6). Let us first consider simple case where $\mathbf{A} = \mathbf{I}_m$ and $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_m)$. Then $Q = X_1^2 + \dots + X_m^2$ has chi-square distribution with m degrees of freedom, $Q \sim \chi_m^2$.

Theorem 3.1 *Let $\mathbf{X} \sim \mathcal{N}_m(\mathbf{0}, \boldsymbol{\Sigma})$. Then $\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X} \sim \chi_m^2$.*

Proof. Consider spectral decomposition $\boldsymbol{\Sigma} = \mathbf{T}\boldsymbol{\Lambda}\mathbf{T}'$ of the covariance matrix $\boldsymbol{\Sigma}$, and random vector $\mathbf{Y} = \boldsymbol{\Sigma}^{-1/2}\mathbf{X}$, where $\boldsymbol{\Sigma}^{-1/2} = \mathbf{T}\boldsymbol{\Lambda}^{-1/2}\mathbf{T}'$. Note that $\mathbb{E}[\mathbf{Y}] = \mathbf{0}$, the covariance matrix of \mathbf{Y} is $\boldsymbol{\Sigma}^{-1/2}\boldsymbol{\Sigma}\boldsymbol{\Sigma}^{-1/2} = \mathbf{I}_m$ and $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_m)$. Moreover

$$\mathbf{Y}'\mathbf{Y} = \mathbf{X}'\boldsymbol{\Sigma}^{-1/2}\boldsymbol{\Sigma}^{-1/2}\mathbf{X} = \mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X}.$$

Hence $\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X} = Y_1^2 + \dots + Y_m^2 \sim \chi_m^2$. □

An $m \times m$ matrix \mathbf{P} is said to be *idempotent* or *projection* matrix if $\mathbf{P}^2 = \mathbf{P}$. All eigenvalues of a projection matrix are either 1 or 0. Indeed if λ is an eigenvalue of \mathbf{P} and \mathbf{e} the corresponding eigenvector, then $\mathbf{P}^2\mathbf{e} = \mathbf{P}(\mathbf{P}\mathbf{e}) = \lambda^2\mathbf{e}$. On the other hand since $\mathbf{P}^2 = \mathbf{P}$, $\mathbf{P}^2\mathbf{e} = \lambda\mathbf{e}$. It follows that $\lambda^2 = \lambda$, and hence $\lambda = 1$ or $\lambda = 0$.

Moreover, suppose that \mathbf{P} is symmetric. Then for any $\mathbf{x} \in \mathbb{R}^m$,

$$(\mathbf{x} - \mathbf{P}\mathbf{x})'\mathbf{P}\mathbf{x} = \mathbf{x}'\mathbf{P}\mathbf{x} - \mathbf{x}'\mathbf{P}'\mathbf{P}\mathbf{x} = \mathbf{x}'\mathbf{P}\mathbf{x} - \mathbf{x}'\mathbf{P}^2\mathbf{x} = 0.$$

That is, \mathbf{P} makes orthogonal projection of vector \mathbf{x} onto the linear space $\{\mathbf{y} : \mathbf{y} = \mathbf{P}\mathbf{x}, \mathbf{x} \in \mathbb{R}^m\}$.

Also by the spectral decomposition, $\mathbf{P} = \mathbf{T}_1\mathbf{T}_1'$, where \mathbf{T}_1 is the $m \times r$ matrix whose columns are orthonormal eigenvectors corresponding to eigenvalues 1, i.e., $\mathbf{T}_1'\mathbf{T}_1 = \mathbf{I}_r$. Then $\text{rank}(\mathbf{P}) = r = \text{tr}(\mathbf{P})$.

Theorem 3.2 *Let $\mathbf{X} \sim \mathcal{N}_m(\mathbf{0}, \mathbf{I}_m)$ and \mathbf{P} be symmetric projection matrix of rank r . Then $\mathbf{X}'\mathbf{P}\mathbf{X} \sim \chi_r^2$.*

Proof. Consider spectral decomposition $\mathbf{P} = \mathbf{T}_1 \mathbf{T}'_1$. Then $\mathbf{X}' \mathbf{P} \mathbf{X} = \mathbf{X}' \mathbf{T}_1 \mathbf{T}'_1 \mathbf{X} = \mathbf{Z}' \mathbf{Z}$, where $\mathbf{Z} = \mathbf{T}'_1 \mathbf{X}$. We have that the $r \times 1$ vector \mathbf{Z} has normal distribution with zero mean vector and covariance matrix $\mathbf{T}'_1 \mathbf{T}_1 = \mathbf{I}_r$. It follows that $\mathbf{X}' \mathbf{P} \mathbf{X} \sim \chi_r^2$. \square

Noncentral chi square distribution.

Let $\mathbf{X} \sim \mathcal{N}_m(\boldsymbol{\mu}, \mathbf{I}_m)$ and consider $Q = \mathbf{X}' \mathbf{X} = X_1^2 + \dots + X_m^2$. Note that if $\mathbf{Y} = \mathbf{T} \mathbf{X}$, where \mathbf{T} is an orthogonal matrix, then $\mathbf{Y}' \mathbf{Y} = \mathbf{X}' \mathbf{X}$ and $\mathbb{E}[\mathbf{Y}] = \mathbf{T} \boldsymbol{\mu}$, and the covariance matrix of \mathbf{Y} is \mathbf{I}_m . It follows that the distribution of Q depends on $\delta = \mu_1^2 + \dots + \mu_m^2$ rather than individual values of the components of the mean vector $\boldsymbol{\mu}$. Distribution of Q is called noncentral chi square with the noncentrality parameter $\delta = \mu_1^2 + \dots + \mu_m^2$ and m degrees of freedom, denoted $Q \sim \chi_m^2(\delta)$. Similar to Theorems 3.1 and 3.2 it is possible to show the following.

Theorem 3.3 *If $\mathbf{X} \sim \mathcal{N}_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then $\mathbf{X}' \boldsymbol{\Sigma}^{-1} \mathbf{X} \sim \chi_m^2(\delta)$ with the noncentrality parameter $\delta = \boldsymbol{\mu}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$. If $\mathbf{X} \sim \mathcal{N}_m(\boldsymbol{\mu}, \mathbf{I}_m)$ and \mathbf{P} is symmetric projection matrix of rank r , then $\mathbf{X}' \mathbf{P} \mathbf{X} \sim \chi_r^2(\delta)$, where $\delta = \boldsymbol{\mu}' \mathbf{P} \boldsymbol{\mu}$.*

4 Statistical inference of linear models

Consider linear regression model

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i, \quad i = 1, \dots, N. \quad (4.1)$$

Denote by $\mathbf{Y} = (Y_1, \dots, Y_N)'$ vector of responses, $\mathbf{X}_j = (x_{1j}, \dots, x_{Nj})'$, $j = 1, \dots, k$, predictors (regressors), $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)'$ vector of errors and $\mathbf{1}_N = (1, \dots, 1)'$ vector of ones. Then we can write model (4.1) as $\mathbf{Y} = \beta_0 \mathbf{1}_N + \beta_1 \mathbf{X}_1 + \dots + \beta_k \mathbf{X}_k + \boldsymbol{\varepsilon}$, or equivalently in matrix form as

$$\mathbf{Y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (4.2)$$

where $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_k)'$ is vector of parameters and $\mathbf{X} = [\mathbf{1}_N, \mathbf{X}_1, \dots, \mathbf{X}_k]$ is $N \times p$, $p = k + 1$, so called design matrix. Note that the first column of \mathbf{X} is formed by ones. Unless stated otherwise, it will be assumed that \mathbf{X} has full column rank p , i.e., column vectors $\mathbf{1}_N, \mathbf{X}_1, \dots, \mathbf{X}_k$ of the design matrix are linearly independent.

Note that the design matrix \mathbf{X} is assumed to be *deterministic*. This is justified when values x_{ij} of the predictors (regressors) are observed without error. If x_{ij} are modelled as random, the analysis below can be pushed through by conditional arguments.

The Least Squares Estimator (LSE) $\hat{\boldsymbol{\beta}}$ of $\boldsymbol{\beta}$ is solution of the problem

$$\min_{\boldsymbol{\beta}} (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X} \boldsymbol{\beta}). \quad (4.3)$$

By Pythagoras Theorem vector of residuals $\mathbf{e} = \mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}$ is orthogonal to the linear space generated by columns of the design matrix \mathbf{X} . That is $\mathbf{e}' \mathbf{X} = \mathbf{0}$ or equivalently $\mathbf{X}' (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}}) = \mathbf{0}$. It follows that $(\mathbf{X}' \mathbf{X}) \hat{\boldsymbol{\beta}} = \mathbf{X}' \mathbf{Y}$. Since it is assumed that matrix \mathbf{X} has full column rank, the $p \times p$ matrix $\mathbf{X}' \mathbf{X}$ is nonsingular (invertible). Thus the LSE can be written as $\hat{\boldsymbol{\beta}} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}$.

Suppose that $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$. Then (recall that the design matrix \mathbf{X} is assumed to be deterministic)

$$\mathbb{E}[\hat{\boldsymbol{\beta}}] = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbb{E}[\mathbf{Y}] = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbb{E}[\mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}] = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' (\mathbf{X} \boldsymbol{\beta} + \mathbb{E}[\boldsymbol{\varepsilon}]) = \boldsymbol{\beta}.$$

That is, $\hat{\boldsymbol{\beta}}$ is an unbiased estimator of $\boldsymbol{\beta}$.

Consider the $N \times N$ matrix $\mathbf{H} := \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'$. Note that vector $\hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\beta}}$ of fitted values is given by $\hat{\mathbf{Y}} = \mathbf{H} \mathbf{Y}$, and vector of residuals $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}}$ is given by $\mathbf{e} = (\mathbf{I}_N - \mathbf{H}) \mathbf{Y}$. Matrix \mathbf{H}

is the orthogonal projection matrix onto the space generated by columns of \mathbf{X} , i.e., $\mathbf{H}\mathbf{Y} = \mathbf{X}\hat{\boldsymbol{\beta}}$ and $\mathbf{X}'(\mathbf{Y} - \mathbf{H}\mathbf{Y}) = \mathbf{0}$; and matrix $\mathbf{I}_N - \mathbf{H}$ is the orthogonal projection matrix onto the space orthogonal to the space generated by columns of \mathbf{X} .

Matrix \mathbf{H} has the following properties:

- (i) \mathbf{H} is symmetric.
- (ii) \mathbf{H} and $\mathbf{I}_N - \mathbf{H}$ are idempotent (projection) matrices, i.e. $\mathbf{H}^2 = \mathbf{H}$ and $(\mathbf{I}_n - \mathbf{H})^2 = \mathbf{I}_n - \mathbf{H}$.
- (iii) $\text{tr}(\mathbf{H}) = p$ and $\text{tr}(\mathbf{I}_N - \mathbf{H}) = N - p$.
- (iv) $\mathbf{H}\mathbf{X} = \mathbf{X}$ and $(\mathbf{I}_N - \mathbf{H})\mathbf{X} = \mathbf{0}$.

Suppose that the errors ε_i , are uncorrelated, $\mathbb{E}(\varepsilon_i) = 0$ and $\text{Var}(\varepsilon_i) = \sigma^2$, $i = 1, \dots, N$, that is, $\text{Cov}(\mathbf{Y}) = \text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2\mathbf{I}_N$. Then the covariance matrix of $\hat{\boldsymbol{\beta}}$ can be computed as

$$\text{Cov}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'[\text{Cov}(\mathbf{Y})]\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}.$$

It also follows that the covariance matrix of \mathbf{e} is

$$\text{Cov}(\mathbf{e}) = \text{Cov}[(\mathbf{I}_N - \mathbf{H})\mathbf{Y}] = \sigma^2(\mathbf{I}_N - \mathbf{H})^2 = \sigma^2(\mathbf{I}_N - \mathbf{H}).$$

Moreover $\mathbb{E}[\mathbf{e}] = (\mathbf{I}_N - \mathbf{H})\mathbb{E}[\mathbf{Y}] = (\mathbf{I}_N - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$ and hence

$$\mathbb{E}[e_1^2 + \dots + e_N^2] = \sum_{i=1}^N \mathbb{E}[e_i^2] = \sum_{i=1}^N \text{Var}(e_i) = \sigma^2 \text{tr}(\mathbf{I}_N - \mathbf{H}) = \sigma^2(N - p).$$

That is,

$$S^2 := \frac{1}{N-p} \sum_{i=1}^N e_i^2$$

is an unbiased estimator of σ^2 .

Since the first column of \mathbf{X} is vector $\mathbf{1}_N = (1, \dots, 1)'$ of ones and $\mathbf{e}'\mathbf{X} = \mathbf{0}$, it follows that $\mathbf{e}'\mathbf{1}_N = 0$, that is $\sum_{i=1}^N e_i = 0$. In a similar way we have $\mathbf{e}'\hat{\mathbf{Y}} = \mathbf{Y}'(\mathbf{I}_N - \mathbf{H})\mathbf{H}\mathbf{Y} = 0$. That is, residuals \mathbf{e} and fitted values $\hat{\mathbf{Y}}$ are uncorrelated.

Consider $\bar{Y} = N^{-1} \sum_{i=1}^N Y_i$. Since $\sum_{i=1}^N e_i = 0$ we have that $\bar{Y} = N^{-1} \sum_{i=1}^N \hat{Y}_i$ as well. Note that

$$\sum_{i=1}^N (\hat{Y}_i - \bar{Y})(Y_i - \hat{Y}_i) = \sum_{i=1}^N (\hat{Y}_i - \bar{Y})e_i = \mathbf{e}'\hat{\mathbf{Y}} - \bar{Y} \sum_{i=1}^N e_i = 0,$$

and hence

$$\sum_{i=1}^N (Y_i - \bar{Y})^2 = \sum_{i=1}^N (Y_i - \hat{Y}_i + \hat{Y}_i - \bar{Y})^2 = \sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^N (Y_i - \hat{Y}_i)^2.$$

That is

$$S_{YY} = SS_R + SS_E, \tag{4.4}$$

where

$$S_{YY} := \sum_{i=1}^N (Y_i - \bar{Y})^2, \quad SS_R := \sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2, \quad SS_E := \sum_{i=1}^N e_i^2. \tag{4.5}$$

The so-called coefficient of determination is defined as

$$R^2 := \frac{S_{YY}}{SS_R} = 1 - \frac{SS_E}{SS_R}.$$

It is interpreted as proportion of the total variation S_{YY} (corrected for the average) explained by variation SS_R due to regression. Another interpretation is that $R^2 = r^2$, where r is the sample correlation coefficient between Y_i and \hat{Y}_i . Indeed

$$\sum_{i=1}^N (\hat{Y}_i - \bar{Y})(Y_i - \bar{Y}) = \sum_{i=1}^N (\hat{Y}_i - \bar{Y})(Y_i - \hat{Y}_i + \hat{Y}_i - \bar{Y}) = \sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2,$$

and hence

$$r^2 = \left[\frac{\sum_{i=1}^N (\hat{Y}_i - \bar{Y})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2} \sqrt{\sum_{i=1}^N (Y_i - \bar{Y})^2}} \right]^2 = \frac{\left[\sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2 \right]^2}{\sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2 \sum_{i=1}^N (Y_i - \bar{Y})^2} = \frac{\sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^N (Y_i - \bar{Y})^2}.$$

In case of one predictor, i.e., $Y_i = \beta_0 + \beta X_i + \varepsilon_i$, $i = 1, \dots, N$, the sample correlation between Y_i and $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta} X_i$ is the same as the sample correlation between Y_i and X_i , $i = 1, \dots, N$.

Theorem 4.1 (Gauss - Markov) *Suppose that $\mathbb{E}[\varepsilon] = \mathbf{0}$ and $\text{Cov}[\varepsilon] = \sigma^2 \mathbf{I}_N$. Then the LSE $\hat{\beta}$ is the Best Linear Unbiased Estimator (BLUE) of β . That is, if $\tilde{\beta} = \mathbf{A}'\mathbf{Y}$ is a linear unbiased estimator of β (i.e., $\mathbb{E}[\tilde{\beta}] = \beta$ for all β), then*

$$\text{Var}(\mathbf{a}'\tilde{\beta}) \geq \text{Var}(\mathbf{a}'\hat{\beta}) \quad (4.6)$$

for any $p \times 1$ vector \mathbf{a} .

Proof. Since $\mathbb{E}[\tilde{\beta}] = \beta$ for all β , it follows that $\beta = \mathbf{A}'\mathbb{E}[\mathbf{Y}] = \mathbf{A}'\mathbf{X}\beta$. Hence $(\mathbf{I}_p - \mathbf{A}'\mathbf{X})\beta = \mathbf{0}$ for all β , and thus $\mathbf{A}'\mathbf{X} = \mathbf{I}_p$. Consider matrix $\mathbf{B} = \mathbf{A} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$. Note that since $\mathbf{X}'\mathbf{A} = \mathbf{I}_p$, it follows that

$$(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{B} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{A} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}) = \mathbf{0},$$

and hence $\mathbf{X}'\mathbf{B} = \mathbf{0}$.

Now since covariance matrix of \mathbf{Y} is $\sigma^2 \mathbf{I}_N$ it follows that

$$\text{Var}(\mathbf{a}'\tilde{\beta}) = \text{Var}(\mathbf{a}'\mathbf{A}'\mathbf{Y}) = \sigma^2 \mathbf{a}'\mathbf{A}'\mathbf{A}\mathbf{a}.$$

Also since $\mathbf{X}'\mathbf{B} = \mathbf{0}$ we have that

$$\mathbf{A}'\mathbf{A} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} + \mathbf{B}'\mathbf{B} = (\mathbf{X}'\mathbf{X})^{-1} + \mathbf{B}'\mathbf{B}.$$

Hence

$$\text{Var}(\mathbf{a}'\tilde{\beta}) = \sigma^2 \mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{a} + \sigma^2 \mathbf{a}'\mathbf{B}'\mathbf{B}\mathbf{a} = \text{Var}(\mathbf{a}'\hat{\beta}) + \sigma^2 \mathbf{a}'\mathbf{B}'\mathbf{B}\mathbf{a}.$$

It remains to note that $\mathbf{a}'\mathbf{B}'\mathbf{B}\mathbf{a} = (\mathbf{B}\mathbf{a})'\mathbf{B}\mathbf{a} \geq 0$. □

The LSE $\hat{\beta}$ is the solution of the system of linear equations $(\mathbf{X}'\mathbf{X})\beta = \mathbf{X}'\mathbf{Y}$. It can happen that small changes in values of the design matrix \mathbf{X} result in big changes in the solution of that system of equations. In numerical analysis such problems are called ill-conditioned. In regression this is called multicollinearity problem, when columns of the design matrix are 'almost' linearly

dependent. Ill conditioning of a system of linear equations is measured by the so-called condition number (see below).

In regression the multicollinearity problem is measured by the so-called Variance Inflation Factor, VIF_i , which is a measure of collinearity of regressor (predictor) \mathbf{X}_i with the other regressors, $i = 1, \dots, k$. It is defined as $VIF_i := 1/(1 - R_i^2)$, where R_i^2 is the coefficient of determination of regression \mathbf{X}_i on the other regressors. Let $\tilde{\mathbf{X}} = [\tilde{\mathbf{X}}_1, \dots, \tilde{\mathbf{X}}_k]$ be the $N \times k$ normalized design matrix, i.e., the averages are removed from each regressor so the sum of elements of each regressor $\tilde{\mathbf{X}}_i$ is zero, and all diagonal elements of $k \times k$ matrix $\mathbf{R} := \tilde{\mathbf{X}}' \tilde{\mathbf{X}}$ are equal to one, i.e., the sum of squared elements of each regressor $\tilde{\mathbf{X}}_i$ is equal to one.

Consider the following partitioning $\begin{bmatrix} 1 & \mathbf{r}_{12} \\ \mathbf{r}_{21} & \mathbf{R}_{11} \end{bmatrix}$ of matrix \mathbf{R} . Then by equation (2.7) we have that the first diagonal element of matrix \mathbf{R}^{-1} is equal to $(1 - \mathbf{r}_{12} \mathbf{R}_{11}^{-1} \mathbf{r}_{21})^{-1}$. Consider now regression of the first regressor $\tilde{\mathbf{X}}_1$ on the other regressors $\tilde{\mathbf{X}}_2, \dots, \tilde{\mathbf{X}}_k$. The corresponding coefficient of determination R_1^2 is equal to $\mathbf{r}_{12} \mathbf{R}_{11}^{-1} \mathbf{r}_{21}$. This can be applied to regression of every $\tilde{\mathbf{X}}_i$ on the other regressors. Therefore Variance Inflation Factors can be obtained as diagonal elements of the matrix \mathbf{R}^{-1} .

Condition number.

Consider the system of linear equations $\mathbf{A}\mathbf{x} = \mathbf{b}$, where \mathbf{A} is a nonsingular $n \times n$ matrix (not necessarily symmetric) and \mathbf{b} is an $n \times 1$ nonzero vector. It has solution $\mathbf{x}_0 = \mathbf{A}^{-1}\mathbf{b}$. Consider perturbed system $\mathbf{A}\mathbf{x} = \mathbf{b} + \boldsymbol{\varepsilon}$, where $\boldsymbol{\varepsilon}$ is a “small” vector of errors. This system has solution $\mathbf{x}_\varepsilon = \mathbf{x}_0 + \mathbf{A}^{-1}\boldsymbol{\varepsilon}$. Consider the following ratio of the relative error in the solution to the relative error in \mathbf{b}

$$\frac{\|\mathbf{A}^{-1}\boldsymbol{\varepsilon}\|/\|\mathbf{x}_0\|}{\|\boldsymbol{\varepsilon}\|/\|\mathbf{b}\|} = \frac{\|\mathbf{A}^{-1}\boldsymbol{\varepsilon}\|}{\|\boldsymbol{\varepsilon}\|} \times \frac{\|\mathbf{b}\|}{\|\mathbf{A}^{-1}\mathbf{b}\|}.$$

The following maximum is called the conditional number of \mathbf{A} :

$$\text{cond}(\mathbf{A}) := \max_{\mathbf{b} \neq \mathbf{0}, \boldsymbol{\varepsilon} \neq \mathbf{0}} \frac{\|\mathbf{A}^{-1}\boldsymbol{\varepsilon}\|}{\|\boldsymbol{\varepsilon}\|} \times \frac{\|\mathbf{b}\|}{\|\mathbf{A}^{-1}\mathbf{b}\|} = \left(\max_{\boldsymbol{\varepsilon} \neq \mathbf{0}} \frac{\|\mathbf{A}^{-1}\boldsymbol{\varepsilon}\|}{\|\boldsymbol{\varepsilon}\|} \right) \left(\max_{\mathbf{b} \neq \mathbf{0}} \frac{\|\mathbf{b}\|}{\|\mathbf{A}^{-1}\mathbf{b}\|} \right).$$

Now let $\sigma_{\min}(\mathbf{A}) = \sqrt{\lambda_{\min}(\mathbf{A}'\mathbf{A})}$ and $\sigma_{\max}(\mathbf{A}) = \sqrt{\lambda_{\max}(\mathbf{A}'\mathbf{A})}$ be the minimal and maximal singular values of \mathbf{A} (see section 15.4). Then

$$\max_{\boldsymbol{\varepsilon} \neq \mathbf{0}} \frac{\|\mathbf{A}^{-1}\boldsymbol{\varepsilon}\|}{\|\boldsymbol{\varepsilon}\|} = \max_{\|\boldsymbol{\varepsilon}\|=1} \sqrt{\boldsymbol{\varepsilon}'(\mathbf{A}'\mathbf{A})^{-1}\boldsymbol{\varepsilon}} = \frac{1}{\sigma_{\min}(\mathbf{A})},$$

and

$$\max_{\mathbf{b} \neq \mathbf{0}} \frac{\|\mathbf{b}\|}{\|\mathbf{A}^{-1}\mathbf{b}\|} = \max_{\mathbf{z} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{z}\|}{\|\mathbf{z}\|} = \sigma_{\max}(\mathbf{A}).$$

Therefore $\text{cond}(\mathbf{A}) = \sigma_{\max}(\mathbf{A})/\sigma_{\min}(\mathbf{A})$.

If matrix \mathbf{A} is symmetric positive definite, then $\sigma_{\max}(\mathbf{A})$ and $\sigma_{\min}(\mathbf{A})$ are the largest and smallest eigenvalues of \mathbf{A} , respectively. For matrix $\mathbf{A} = \gamma\mathbf{I}_n$, $\gamma \neq 0$, its condition number $\text{cond}(\gamma\mathbf{I}_n) = 1$. Otherwise the condition number is bigger than one.

4.1 Distribution theory

Suppose now that $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I}_N)$ and hence $\mathbf{Y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I}_N)$. It follows that the LSE $\hat{\boldsymbol{\beta}}$ has normal distribution $\mathcal{N}(\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1})$. Hence it follows by Theorem 3.1 that

$$\sigma^{-2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{X}' \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sim \chi_p^2, \quad (4.7)$$

where $p = k + 1$. Recall that $S^2 = (N - p)^{-1} \mathbf{e}' \mathbf{e}$ is an unbiased estimator of σ^2 .

Since $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ and $(\mathbf{I}_N - \mathbf{H})\mathbf{X} = \mathbf{0}$,

$$\frac{(N - p)S^2}{\sigma^2} = \frac{\mathbf{e}' \mathbf{e}}{\sigma^2} = \frac{\mathbf{Y}'(\mathbf{I}_N - \mathbf{H})^2 \mathbf{Y}}{\sigma^2} = \frac{\mathbf{Y}'(\mathbf{I}_N - \mathbf{H})\mathbf{Y}}{\sigma^2} = \frac{\boldsymbol{\varepsilon}'(\mathbf{I}_N - \mathbf{H})\boldsymbol{\varepsilon}}{\sigma^2}.$$

Recall that $\mathbf{I}_N - \mathbf{H}$ is a projection matrix. Its rank

$$\begin{aligned} \text{rank}(\mathbf{I}_N - \mathbf{H}) &= \text{tr}(\mathbf{I}_N - \mathbf{H}) = \text{tr}(\mathbf{I}_N) - \text{tr}(\mathbf{H}) = N - \text{tr}(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') \\ &= N - \text{tr}((\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}) = N - p. \end{aligned}$$

By Theorem 3.2 it follows that

$$\frac{(N - p)S^2}{\sigma^2} \sim \chi_{N-p}^2. \quad (4.8)$$

Moreover

$$\text{Cov}[\mathbf{e}, \hat{\boldsymbol{\beta}}] = (\mathbf{I}_N - \mathbf{H})\text{Cov}(\mathbf{Y})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{I}_N - \mathbf{H})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \mathbf{0}.$$

Hence \mathbf{e} and $\hat{\boldsymbol{\beta}}$ are independent. It follows that S^2 and $\hat{\boldsymbol{\beta}}$ are independent, and hence S^2 and $\sigma^{-2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})'\mathbf{X}'\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$ are independent. It follows that

$$\frac{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})'\mathbf{X}'\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})/p}{S^2} \sim F_{p, N-p}. \quad (4.9)$$

This can be used to construct the following $(1 - \alpha)$ -confidence region for $\boldsymbol{\beta}$:

$$\{\boldsymbol{\beta} : (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})'\mathbf{X}'\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \leq pS^2 F_{\alpha; p, N-p}\}.$$

Now consider S_{YY} , SS_R and SS_E and recall that $S_{YY} = SS_R + SS_E$ (see equation (4.4)).

Also

$$(N - p)S^2 = SS_E = \mathbf{Y}'(\mathbf{I}_N - \mathbf{H})\mathbf{Y} \quad (4.10)$$

and

$$SS_R = (\hat{\mathbf{Y}} - \mathbf{1}_N \bar{Y})'(\hat{\mathbf{Y}} - \mathbf{1}_N \bar{Y}) = (\mathbf{H}\mathbf{Y} - N^{-1}\mathbf{1}_N \mathbf{1}'_N \mathbf{Y})'(\mathbf{H}\mathbf{Y} - N^{-1}\mathbf{1}_N \mathbf{1}'_N \mathbf{Y}) = \mathbf{Y}'(\mathbf{H} - N^{-1}\mathbf{1}_N \mathbf{1}'_N)^2 \mathbf{Y}.$$

Moreover, since $\mathbf{H}\mathbf{1}_N = \mathbf{1}_N$ (this holds since $\mathbf{H}\mathbf{X} = \mathbf{X}$ and the first column of \mathbf{X} is $\mathbf{1}_N$) we obtain $(\mathbf{H} - N^{-1}\mathbf{1}_N \mathbf{1}'_N)^2 = \mathbf{H} - N^{-1}\mathbf{1}_N \mathbf{1}'_N$. and hence

$$SS_R = \mathbf{Y}'(\mathbf{H} - N^{-1}\mathbf{1}_N \mathbf{1}'_N)\mathbf{Y}. \quad (4.11)$$

Since $(\mathbf{I}_N - \mathbf{H})\mathbf{H} = \mathbf{0}$ and $(\mathbf{I}_N - \mathbf{H})\mathbf{1}_N = \mathbf{0}$, we have that

$$(\mathbf{I}_N - \mathbf{H})(\mathbf{H} - N^{-1}\mathbf{1}_N \mathbf{1}'_N) = (\mathbf{I}_N - \mathbf{H})\mathbf{H} - N^{-1}(\mathbf{I}_N - \mathbf{H})\mathbf{1}_N \mathbf{1}'_N = \mathbf{0},$$

and hence SS_E and SS_R are independent.

Consider the following so-called F -statistic, for testing $H_0 : \beta_1 = \dots = \beta_k = 0$, against the alternative that at least one $\beta_i \neq 0$,

$$F = \frac{SS_R/k}{SS_E/(N - p)}. \quad (4.12)$$

Recall that $SS_E/\sigma^2 \sim \chi_{N-p}^2$. Also under H_0 we have that $\mathbf{Y} = \beta_0 \mathbf{1}_N$ and hence

$$(\mathbf{H} - N^{-1}\mathbf{1}_N \mathbf{1}'_N)\mathbf{Y} = \beta_0(\mathbf{H} - N^{-1}\mathbf{1}_N \mathbf{1}'_N)\mathbf{1}_N = \beta_0(\mathbf{H}\mathbf{1}_N - N^{-1}\mathbf{1}_N \mathbf{1}'_N \mathbf{1}_N) = \mathbf{1}_N - \mathbf{1}_N = \mathbf{0}.$$

Consequently

$$SS_R = \boldsymbol{\varepsilon}'(\mathbf{H} - N^{-1}\mathbf{1}_N\mathbf{1}'_N)\boldsymbol{\varepsilon},$$

and hence $SS_R/\sigma^2 \sim \chi_k^2$. Note that

$$\text{rank}(\mathbf{H} - N^{-1}\mathbf{1}_N\mathbf{1}'_N) = \text{tr}(\mathbf{H} - N^{-1}\mathbf{1}_N\mathbf{1}'_N) = \text{tr}(\mathbf{H}) - 1 = k.$$

It follows that under H_0 the statistic F has $F_{k, N-p}$ distribution.

Under alternative H_1 , SS_R/σ^2 has noncentral chi square distribution $SS_R/\sigma^2 \sim \chi_k^2(\delta)$ with noncentrality parameter

$$\delta = \sigma^{-2}\boldsymbol{\beta}'\mathbf{X}'(\mathbf{H} - N^{-1}\mathbf{1}_N\mathbf{1}'_N)\mathbf{X}\boldsymbol{\beta} = \sigma^{-2}\boldsymbol{\beta}'(\mathbf{X}'\mathbf{X} - N^{-1}(\mathbf{1}'_N\mathbf{X})'(\mathbf{1}'_N\mathbf{X}))\boldsymbol{\beta}.$$

Therefore under the alternative, the F -statistic has noncentral F distribution with the noncentrality parameter δ for SS_R .

4.2 Estimation with linear constraints

Suppose that we want to test linear constraints $\mathbf{a}'_i\boldsymbol{\beta} = c_i$, $i = 1, \dots, q$. We can write this as $\mathbf{A}\boldsymbol{\beta} = \mathbf{c}$, where \mathbf{A} is the corresponding $q \times p$ matrix whose rows are formed from vectors \mathbf{a}'_i , $i = 1, \dots, q$, and $\mathbf{c} = (c_1, \dots, c_q)'$. We assume that vectors \mathbf{a}_i , $i = 1, \dots, q$, are linearly independent, i.e., matrix \mathbf{A} has full row rank q .

The respective constrained least squares estimator $\hat{\boldsymbol{\beta}}_H$ is obtained as a solution of the following optimization problem

$$\min_{\boldsymbol{\beta}} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \text{ subject to } \mathbf{A}\boldsymbol{\beta} = \mathbf{c}. \quad (4.13)$$

Consider the Lagrangian of the above problem (4.13):

$$\begin{aligned} L(\boldsymbol{\beta}, \boldsymbol{\lambda}) &:= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + 2 \sum_{i=1}^q \lambda_i (\mathbf{a}'_i\boldsymbol{\beta} - c_i) \\ &= (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + 2\boldsymbol{\lambda}'(\mathbf{A}\boldsymbol{\beta} - \mathbf{c}). \end{aligned}$$

Problem (4.13) is a convex quadratic problem. Optimality conditions for problem (4.13) can be written as $\partial L(\boldsymbol{\beta}, \boldsymbol{\lambda})/\partial \boldsymbol{\beta} = \mathbf{0}$ and $\mathbf{A}\boldsymbol{\beta} = \mathbf{c}$. Note that

$$\partial L(\boldsymbol{\beta}, \boldsymbol{\lambda})/\partial \boldsymbol{\beta} = -2\mathbf{X}'\mathbf{Y} + 2\mathbf{X}'\mathbf{X}\boldsymbol{\beta} + 2\mathbf{A}'\boldsymbol{\lambda}.$$

Hence the optimality conditions can be written as the following system of linear equations

$$\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{A}' \\ \mathbf{A} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta} \\ \boldsymbol{\lambda} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{Y} \\ \mathbf{c} \end{bmatrix}. \quad (4.14)$$

Note that the Schur complement of matrix $\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{A}' \\ \mathbf{A} & \mathbf{0} \end{bmatrix}$ with respect to matrix $\mathbf{X}'\mathbf{X}$ is $-\mathbf{A}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}'$. Since it is assumed that matrix \mathbf{X} has full column rank and hence matrix $\mathbf{X}'\mathbf{X}$ is nonsingular, and since matrix \mathbf{A} has full row rank q , it follows that $-\mathbf{A}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}'$ is nonsingular, and hence matrix $\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{A}' \\ \mathbf{A} & \mathbf{0} \end{bmatrix}$ is invertible. Therefore the corresponding estimators are given by

$$\begin{bmatrix} \hat{\boldsymbol{\beta}}_H \\ \hat{\boldsymbol{\lambda}}_H \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{A}' \\ \mathbf{A} & \mathbf{0} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{X}'\mathbf{Y} \\ \mathbf{c} \end{bmatrix}. \quad (4.15)$$

By using formula (2.7) for the inverse $\begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{A}' \\ \mathbf{A} & \mathbf{0} \end{bmatrix}^{-1}$, after some algebraic calculations it is possible to write the estimator $\hat{\boldsymbol{\beta}}_H$ in the following form

$$\hat{\boldsymbol{\beta}}_H = \hat{\boldsymbol{\beta}} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}'[\mathbf{A}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}']^{-1}(\mathbf{c} - \mathbf{A}\hat{\boldsymbol{\beta}}), \quad (4.16)$$

where $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ is the unconstrained LSE. It is possible to give the following geometrical interpretation. Consider $\hat{\mathbf{Y}}_H = \mathbf{X}\hat{\boldsymbol{\beta}}_H$. Recall that $\mathbf{Y} - \hat{\mathbf{Y}}$ is orthogonal to the linear space generated by columns of matrix \mathbf{X} . Since $\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_H = \mathbf{X}(\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_H)$, it follows that $\mathbf{Y} - \hat{\mathbf{Y}}$ is orthogonal to $\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_H$. Hence (Pythagoras Theorem)

$$\|\mathbf{Y} - \hat{\mathbf{Y}}_H\|^2 = \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 + \|\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_H\|^2, \quad (4.17)$$

where $\|\cdot\|$ is the Euclidean norm. Moreover

$$\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_H = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}'[\mathbf{A}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}']^{-1}(\mathbf{c} - \mathbf{A}\hat{\boldsymbol{\beta}}). \quad (4.18)$$

The term $\|\mathbf{Y} - \hat{\mathbf{Y}}_H\|^2$ represents the sum of squares of residuals of the reduced (constrained) model, i.e., it is the optimal value of the least squares problem (4.13), and the term $\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2$ is the sum of squares of residuals of the full (unconstrained) model. By (4.18),

$$\|\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_H\|^2 = (\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{c})'[\mathbf{A}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}']^{-1}(\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{c}). \quad (4.19)$$

The F -statistic for testing $H_0 : \mathbf{A}\boldsymbol{\beta} = \mathbf{c}$ is

$$F = \frac{(SS_E(H) - SS_E(F))/q}{SS_E(F)/(N - p)}, \quad (4.20)$$

where $SS_E(F)$ is the sum of squares of residuals of the full (unconstrained) model and $SS_E(H)$ is the sum of squares of residuals of the reduced (constrained) model. By (4.17) and (4.19) we have

$$SS_E(H) - SS_E(F) = (\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{c})'[\mathbf{A}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}']^{-1}(\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{c}). \quad (4.21)$$

Recall that $\hat{\boldsymbol{\beta}} \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1})$, and hence

$$\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{c} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\beta} - \mathbf{c}, \sigma^2\mathbf{A}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}').$$

It follows by Theorem 3.1 that under the $H_0 : \mathbf{A}\boldsymbol{\beta} = \mathbf{c}$ hypothesis, $[SS_E(H) - SS_E(F)]/\sigma^2 \sim \chi_q^2$. Also $SS_E(F)/\sigma^2 \sim \chi_{N-p}^2$ and $SS_E(F)$ is independent of $\hat{\boldsymbol{\beta}}$, and hence $SS_E(H) - SS_E(F)$ and $SS_E(F)$ are independent. It follows that under the H_0 hypothesis, the F statistic (4.20) has $F_{q, N-p}$ distribution.

The F -statistic (4.12), for testing $H_0 : \beta_1 = \dots = \beta_k = 0$, is a particular case of the F -statistic (4.20). Indeed in that case, under H_0 , the LSE $\hat{\boldsymbol{\beta}}_0 = \bar{\mathbf{Y}}$ and hence $SS_E(H) = S_{YY}$. It follows that $SS_E(H) - SS_E(F) = SS_R$.

The statistical inference discussed in section 4.1 and this section is based on the assumption that the error vector $\boldsymbol{\varepsilon}$ has normal distribution. Without this assumption the inference is asymptotic. Recall that for large N , the $qF_{q, N-p}$ distribution becomes approximately like χ_q^2 distribution.

4.3 Polynomial Regression

Consider the polynomial regression model (one predictor)

$$Y_i = \beta_0 + \beta_1 x_i + \dots + \beta_k x_i^k + \varepsilon_i, \quad i = 1, \dots, N. \quad (4.22)$$

We can formulate this as the linear multivariate model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ with the design matrix

$$\mathbf{X} = \begin{bmatrix} 1 & x_1 & \cdots & x_1^k \\ \cdots & \cdots & \cdots & \cdots \\ 1 & x_N & \cdots & x_N^k \end{bmatrix}.$$

We have here $[\mathbf{X}'\mathbf{X}]_{st} = \sum_{i=1}^N x_i^{s+t}$.

Note that (Riemann sum)

$$\int_0^1 x^{s+t} dx \approx \frac{1}{N} \sum_{i=1}^N x_i^{s+t},$$

where x_i is a point of the interval $[(i-1)/N, i/N]$, $i = 1, \dots, N$. Therefore

$$[\mathbf{X}'\mathbf{X}]_{st} = \sum_{i=1}^N x_i^{s+t} \approx N \int_0^1 x^{s+t} dx = \frac{N}{s+t+1}, \quad s, t = 0, \dots, k.$$

That is

$$\mathbf{X}'\mathbf{X} \approx N \begin{bmatrix} 1 & 1/2 & 1/3 & \cdots & 1/(k+1) \\ 1/2 & 1/3 & 1/4 & \cdots & 1/(k+2) \\ 1/3 & 1/4 & 1/5 & \cdots & 1/(k+3) \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 1/(k+1) & 1/(k+2) & 1/(k+3) & \cdots & 1/(2k+1) \end{bmatrix}.$$

This matrix is ill conditioned. Therefore polynomial regression of the form (4.22) typically has multicollinearity problem for $k \geq 3$. To a certain extent this can be dealt with by using orthogonal polynomials. A famous example of orthogonal polynomials is Chebishev polynomials. Even so, polynomial regression of degree larger than 2 usually is difficult to interpret.

Chebishev polynomials

$$T_m(x) = \cos[m(\arccos x)], \quad -1 \leq x \leq 1.$$

Let $\theta = \arccos x$. Then

$$\begin{aligned} T_0(x) &= \cos 0 = 1, \\ T_1(x) &= \cos \theta = x, \\ T_2(x) &= \cos(2\theta) = 2 \cos^2 \theta - 1 = 2x^2 - 1. \end{aligned}$$

Recall that

$$\cos(m+1)\theta + \cos(m-1)\theta = 2 \cos \theta \cos m\theta.$$

It follows that

$$T_{m+1}(x) + T_{m-1}(x) = 2xT_m(x),$$

and hence

$$T_{m+1}(x) = 2xT_m(x) - T_{m-1}(x)$$

can be used for recursive computation of Chebishev polynomials. For example

$$T_3(x) = 2xT_2(x) - T_1(x) = 2x(2x^2 - 1) - x.$$

By using substitution $d \arccos x = \frac{1}{\sqrt{1-x^2}} dx$, we can compute the following integral

$$\int_{-1}^1 \frac{T_k(x)T_\ell(x)}{\sqrt{1-x^2}} dx = \int_{-1}^1 \cos(k\theta) \cos(\ell\theta) d\theta = 0, \text{ for } k \neq \ell.$$

For $x_i = \cos(\pi/N)i$ and $\theta_i = (\pi/N)i$ we have

$$\sum_{i=0}^{N-1} T_k(x_i)T_\ell(x_i) = \sum_{i=0}^{N-1} \cos k\theta_i \cos \ell\theta_i = 0, \quad k \neq \ell. \quad (4.23)$$

For the corresponding polynomial regression

$$Y_i = \beta_0 T_0(x_i) + \beta_1 T_1(x_i) + \dots + \beta_k T_k(x_i) + \varepsilon_i, \quad i = 1, \dots, N,$$

the design matrix is

$$\mathbf{X} = \begin{bmatrix} T_0(x_1) & \cdots & T_k(x_1) \\ \vdots & \ddots & \vdots \\ T_0(x_N) & \cdots & T_k(x_N) \end{bmatrix}.$$

By (4.23) columns of this design matrix are orthogonal to each other, and hence matrix $\mathbf{X}'\mathbf{X}$ is diagonal.

5 Shrinkage Methods

A norm $\|\cdot\|$, on space \mathbb{R}^m , assigns a nonnegative number to vector $\mathbf{x} \in \mathbb{R}^m$. It should have the following properties: (i) $\|\mathbf{x}\| > 0$ for any $\mathbf{x} \neq \mathbf{0}$, (ii) $\|\lambda\mathbf{x}\| = |\lambda| \|\mathbf{x}\|$ for any $\lambda \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^m$, (iii) $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^m$. Properties (ii) and (iii) imply that function $f(\mathbf{x}) = \|\mathbf{x}\|$ is convex. Any two norms $\|\cdot\|$ and $\|\cdot\|'$ on \mathbb{R}^m are equivalent in the following sense: there is a constant $C > 0$ (depending on dimension m of the space \mathbb{R}^m) such that $\|\mathbf{x}\| \leq C\|\mathbf{x}\|'$ and $\|\mathbf{x}\|' \leq C\|\mathbf{x}\|$ for all $\mathbf{x} \in \mathbb{R}^m$.

Important examples of norms are the ℓ_q , $q \geq 1$, norms defined as $\|\mathbf{x}\|_q = (|x_1|^q + \dots + |x_m|^q)^{1/q}$. In particular, the ℓ_2 norm is the Euclidean norm $\|\mathbf{x}\|_2 = \sqrt{x_1^2 + \dots + x_m^2}$, and ℓ_1 norm is $\|\mathbf{x}\|_1 = |x_1| + \dots + |x_m|$. Note that function $\|\cdot\|_q$ is homogeneous (i.e., satisfies the above property (ii)) for any $q > 0$. However for $q \in (0, 1)$, $\|\cdot\|_q$ does not satisfy property (iii), i.e., it is not convex.

5.1 Ridge Regression

Consider the following approach, called Ridge Regression, to estimation parameters of the linear model (4.2)

$$\min_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \varepsilon\|\beta\|_2^2, \quad (5.1)$$

where $\varepsilon > 0$. Solution $\tilde{\beta}_\varepsilon$ of this problem satisfies optimality conditions

$$-\mathbf{X}'(\mathbf{Y} - \mathbf{X}\beta) + \varepsilon\beta = \mathbf{0}.$$

That is $\tilde{\beta}_\varepsilon = (\mathbf{X}'\mathbf{X} + \varepsilon\mathbf{I}_p)^{-1}\mathbf{X}'\mathbf{Y}$ (recall that $p = k + 1$ is the number of estimated parameters). Of course for $\varepsilon = 0$ the estimator $\tilde{\beta}_\varepsilon$ coincides with the LSE $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$. It is also possible to formulate problem (5.1) in the following form

$$\min_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 \quad \text{subject to } \|\beta\|_2 \leq c, \quad (5.2)$$

for a certain value of $c > 0$ (take $c = \|\tilde{\beta}_\varepsilon\|_2$). Conversely solution of problem (5.2), for some $c > 0$, is also the solution of problem (5.1) when ε is the corresponding Lagrange multiplier. (If $\|\hat{\beta}\|_2 \leq c$, then the corresponding $\varepsilon = 0$.) Therefore in a sense problems (5.1) and (5.2) are equivalent to each other for a proper choice of the respective positive constants ε and c .

The estimator $\tilde{\beta}_\varepsilon$ shrinks the LSE to the origin. In particular if columns of the design matrix \mathbf{X} are orthogonal, i.e., matrix $\mathbf{X}'\mathbf{X} = \text{diag}(\lambda_1, \dots, \lambda_p)$ is diagonal. Then

$$\mathbf{X}'\mathbf{X} + \varepsilon\mathbf{I}_p = \text{diag}(\lambda_1 + \varepsilon, \dots, \lambda_p + \varepsilon)$$

and $\tilde{\beta}_{\varepsilon,i} = (1 + \varepsilon/\lambda_i)^{-1}\hat{\beta}_i$. Let $\mathbf{X}'\mathbf{X} = \mathbf{T}\mathbf{\Lambda}\mathbf{T}'$ be the spectral decomposition of matrix $\mathbf{X}'\mathbf{X}$, with $\lambda_1 \geq \dots \geq \lambda_p > 0$ being the eigenvalues and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_p)$. Then $\mathbf{X}'\mathbf{X} + \varepsilon\mathbf{I}_p = \mathbf{T}(\mathbf{\Lambda} + \varepsilon\mathbf{I}_p)\mathbf{T}'$.

Recall that number λ_1/λ_p is called the condition number of matrix $\mathbf{X}'\mathbf{X}$. The condition number of matrix $\mathbf{X}'\mathbf{X} + \varepsilon\mathbf{I}_p$ is $(\lambda_1 + \varepsilon)/(\lambda_p + \varepsilon)$, and can be much smaller than λ_1/λ_p even for small values of $\varepsilon > 0$ if the ratio ε/λ_p is large. Moreover $\tilde{\beta}_\varepsilon = \mathbf{T}(\mathbf{\Lambda} + \varepsilon\mathbf{I}_p)^{-1}\mathbf{T}'\mathbf{X}'\mathbf{Y}$, and hence

$$\tilde{\gamma}_\varepsilon = (\mathbf{\Lambda} + \varepsilon\mathbf{I}_p)^{-1}\tilde{\mathbf{X}}'\mathbf{Y},$$

where $\tilde{\gamma}_\varepsilon = \mathbf{T}'\tilde{\beta}_\varepsilon$ and $\tilde{\mathbf{X}} = \mathbf{X}\mathbf{T}$. Note that $\tilde{\mathbf{X}}'\tilde{\mathbf{X}} = \mathbf{\Lambda}$ and hence $\tilde{\gamma}_{\varepsilon,i} = (1 + \varepsilon/\lambda_i)^{-1}\hat{\gamma}_i$, where $\hat{\gamma}$ is the LSE of the corresponding linear model with \mathbf{X} replaces by $\tilde{\mathbf{X}}$. If ε is much larger than λ_i , and hence the ratio ε/λ_i is large, then $\tilde{\gamma}_{\varepsilon,i}$ becomes small. In that sense this procedure removes from the design matrix $\tilde{\mathbf{X}}$ columns corresponding to small values of the eigenvalues λ_i , and in an implicit way is related to the Principal Components Analysis discussed in section 15.

The estimator $\tilde{\beta}_\varepsilon$ is biased, that is $\mathbb{E}[\tilde{\beta}_\varepsilon] = (\mathbf{X}'\mathbf{X} + \varepsilon\mathbf{I}_p)^{-1}\mathbf{X}'\mathbf{X}\beta$. It is possible to show that there exists $\varepsilon > 0$ such that the components of $\tilde{\beta}_\varepsilon$ have smaller Mean Square Error (MSE) than the respective components of the LSE $\hat{\beta}$. That is, let $\theta = \mathbf{a}'\beta$ for some given vector $\mathbf{a} \neq \mathbf{0}$, and let $\tilde{\theta}_\varepsilon = \mathbf{a}'\tilde{\beta}_\varepsilon$ and $\hat{\theta} = \mathbf{a}'\hat{\beta}$ be estimators of θ . Note that $\mathbf{a}'\hat{\beta}$ is an unbiased estimator of $\mathbf{a}'\beta$. We show that there exists $\varepsilon > 0$ such that

$$MSE(\tilde{\theta}_\varepsilon) < MSE(\hat{\theta}),$$

where $MSE(\tilde{\theta}) = \mathbb{E}[(\tilde{\theta} - \theta)^2]$ is the mean square error of an estimator $\tilde{\theta}$.

Recall that

$$\tilde{\beta}_\varepsilon = (\mathbf{X}'\mathbf{X} + \varepsilon\mathbf{I}_p)^{-1}\mathbf{X}'\mathbf{Y} = [\mathbf{I}_p + \varepsilon(\mathbf{X}'\mathbf{X})^{-1}]^{-1}\hat{\beta},$$

and hence

$$\mathbb{E}[\tilde{\theta}_\varepsilon] = \mathbf{a}'[\mathbf{I}_p + \varepsilon(\mathbf{X}'\mathbf{X})^{-1}]^{-1}\beta.$$

For a matrix \mathbf{A} sufficiently small we have the following geometric series expansion

$$(\mathbf{I} + \mathbf{A})^{-1} = \mathbf{I} - \mathbf{A} + \mathbf{A}^2 - \dots = \mathbf{I} - \mathbf{A} + o(\|\mathbf{A}\|),$$

where $\|\mathbf{A}\| := \sup_{\mathbf{x} \neq \mathbf{0}} \|\mathbf{A}\mathbf{x}\|/\|\mathbf{x}\|$. By applying this to matrix $\mathbf{A} = \varepsilon(\mathbf{X}'\mathbf{X})^{-1}$ for $\varepsilon > 0$ small enough, we obtain

$$\mathbb{E}[\tilde{\theta}_\varepsilon] = \mathbf{a}'[\mathbf{I}_p - \varepsilon(\mathbf{X}'\mathbf{X})^{-1}]\beta + o(\varepsilon) = \theta - \varepsilon\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\beta + o(\varepsilon),$$

where $o(\varepsilon)/\varepsilon \rightarrow 0$ as $\varepsilon \rightarrow 0$. It follows that

$$\text{Bias}[\tilde{\theta}_\varepsilon] = \mathbb{E}[\tilde{\theta}_\varepsilon] - \theta = -\varepsilon \mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\boldsymbol{\beta} + o(\varepsilon),$$

and hence

$$\text{Bias}^2[\tilde{\theta}_\varepsilon] = \varepsilon^2[\mathbf{a}'(\mathbf{X}'\mathbf{X})^{-1}\boldsymbol{\beta}]^2 + o(\varepsilon^2) = o(\varepsilon).$$

We also have that

$$\begin{aligned} \text{Var}[\tilde{\theta}_\varepsilon] &= \sigma^2 \mathbf{a}'[\mathbf{I}_p + \varepsilon(\mathbf{X}'\mathbf{X})^{-1}]^{-1}(\mathbf{X}\mathbf{X}')^{-1}[\mathbf{I}_p + \varepsilon(\mathbf{X}'\mathbf{X})^{-1}]^{-1}\mathbf{a} \\ &= \sigma^2 \mathbf{a}'[\mathbf{I}_p - \varepsilon(\mathbf{X}'\mathbf{X})^{-1}](\mathbf{X}\mathbf{X}')^{-1}[\mathbf{I}_p - \varepsilon(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{a} + o(\varepsilon) \\ &= \sigma^2 \mathbf{a}'(\mathbf{X}\mathbf{X}')^{-1}\mathbf{a} - 2\varepsilon\sigma^2 \mathbf{a}'(\mathbf{X}'\mathbf{X})^{-2}\mathbf{a} + o(\varepsilon) \\ &= \text{Var}[\hat{\theta}] - 2\varepsilon\sigma^2 \mathbf{a}'(\mathbf{X}'\mathbf{X})^{-2}\mathbf{a} + o(\varepsilon). \end{aligned}$$

Therefore

$$MSE(\hat{\theta}) - MSE(\tilde{\theta}_\varepsilon) = \text{Var}[\hat{\theta}] - \text{Var}[\tilde{\theta}_\varepsilon] - \text{Bias}^2[\tilde{\theta}_\varepsilon] = 2\varepsilon\sigma^2 \mathbf{a}'(\mathbf{X}'\mathbf{X})^{-2}\mathbf{a} + o(\varepsilon).$$

Since matrix $\mathbf{X}'\mathbf{X}$ is positive definite, and hence $(\mathbf{X}'\mathbf{X})^{-2}$ is positive definite, and $\mathbf{a} \neq \mathbf{0}$, we have that $\sigma^2 \mathbf{a}'(\mathbf{X}'\mathbf{X})^{-2}\mathbf{a} > 0$. It follows that for $\varepsilon > 0$ small enough the term $2\varepsilon\sigma^2 \mathbf{a}'(\mathbf{X}'\mathbf{X})^{-2}\mathbf{a} + o(\varepsilon)$ is positive, and hence $MSE(\tilde{\theta}_\varepsilon) < MSE(\hat{\theta})$.

In particular this implies that for every $i \in \{1, \dots, k\}$, there exists $\varepsilon > 0$ such that $MSE(\tilde{\beta}_{\varepsilon,i}) < MSE(\hat{\beta}_i)$. However for different $i \in \{1, \dots, k\}$ the corresponding ε can be different, and could be difficult to find. In practical applications the components of $\tilde{\beta}_\varepsilon$ are plotted as a function of $\varepsilon > 0$ until they stabilize.

5.2 Lasso method

The Least Absolute Shrinkage and Selection Operator (Lasso) method is based on using regularization term of the form $\varepsilon\|\boldsymbol{\beta}\|_1$ for some $\varepsilon > 0$. That is, the Lasso estimator $\tilde{\boldsymbol{\beta}}_\varepsilon$ is obtained as a solution of the following optimization problem

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \varepsilon\|\boldsymbol{\beta}\|_1. \quad (5.3)$$

Equivalently this can be formulated as

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \quad \text{subject to } \|\boldsymbol{\beta}\|_1 \leq c, \quad (5.4)$$

for an appropriate choice of the constant $c > 0$. If $c < \|\hat{\boldsymbol{\beta}}\|_1$, then the Lasso estimator performs shrinkage of the LSE $\hat{\boldsymbol{\beta}}$.

Note that

$$\frac{\partial \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2}{\partial \boldsymbol{\beta}} = 2(\mathbf{X}'\mathbf{X}\boldsymbol{\beta} - \mathbf{X}'\mathbf{Y}) = 2(\mathbf{X}'\mathbf{X})(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}),$$

where $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ is the usual least squares estimator. When $c < \|\hat{\boldsymbol{\beta}}\|_1$, an optimal solution of problem (5.4) is on the boundary of the feasible set $S = \{\boldsymbol{\beta} : \|\boldsymbol{\beta}\|_1 \leq c\}$ and the corresponding optimality conditions are

$$-2(\mathbf{X}'\mathbf{X})(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \in N_S(\boldsymbol{\beta}),$$

where $N_S(\boldsymbol{\beta}) := \{\boldsymbol{\gamma} : \boldsymbol{\gamma}'(\boldsymbol{\zeta} - \boldsymbol{\beta}) \leq 0, \forall \boldsymbol{\zeta} \in S\}$ is the normal cone to S at $\boldsymbol{\beta} \in S$.

Optimality conditions for problem (5.3) are

$$\mathbf{0} \in 2(\mathbf{X}'\mathbf{X})(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \varepsilon\partial\|\boldsymbol{\beta}\|_1,$$

where $\partial\|\boldsymbol{\beta}\|_1$ is the subdifferential of the function $f(\boldsymbol{\beta}) = \|\boldsymbol{\beta}\|_1$. The subdifferential $\partial\|\boldsymbol{\beta}\|_1$ consists of vectors \mathbf{g} (the so-called subgradients) such that $g_i = 1$ if $\beta_i > 0$, $g_i = -1$ if $\beta_i < 0$, and g_i can be any number of the interval $[-1, 1]$ if $\beta_i = 0$. It follows that if ε is bigger than the absolute value of every component $[\mathbf{X}'\mathbf{Y}]_i$ of vector $\mathbf{X}'\mathbf{Y}$, then $\tilde{\boldsymbol{\beta}}_\varepsilon = \mathbf{0}$. If $\mathbf{X}'\mathbf{X}$ is diagonal, then $\tilde{\beta}_{\varepsilon,i} = 0$, when $\varepsilon > [\mathbf{X}'\mathbf{Y}]_i$.

It is possible to look at Lasso estimation from the following point of view. By definition $\|\boldsymbol{\beta}\|_0$ is equal to the number of nonzero components of vector $\boldsymbol{\beta}$. Note that $\|\boldsymbol{\beta}\|_0 = \lim_{q \downarrow 0} \sum |\beta_i|^q$. Consider the problem

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \quad \text{subject to } \|\boldsymbol{\beta}\|_0 \leq c, \quad (5.5)$$

i.e., it is the least squares problem subject to the constraint that the number of used regressors is not larger than c . This is a difficult combinatorial problem. Problem (5.4) can be viewed as a convex approximation of problem (5.5). Problem (5.3) can be formulated as the following problem

$$\begin{aligned} \min_{\boldsymbol{\beta}, \boldsymbol{\xi}} \quad & \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \varepsilon \sum_{i=0}^k \xi_i \\ \text{s.t.} \quad & \beta_i \leq \xi_i, \quad -\beta_i \leq \xi_i, \quad i = 0, \dots, k, \end{aligned} \quad (5.6)$$

and problem (5.4) as

$$\begin{aligned} \min_{\boldsymbol{\beta}, \boldsymbol{\xi}} \quad & \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \\ \text{s.t.} \quad & \beta_i \leq \xi_i, \quad -\beta_i \leq \xi_i, \quad i = 0, \dots, k, \\ & \sum_{i=0}^k \xi_i \leq c. \end{aligned} \quad (5.7)$$

Both problems (5.6) and (5.7) are convex quadratic programming problems, and can be solved efficiently.

6 Elements of large samples theory

Let Y_n , $n = 1, \dots$, be a sequence of random variables. It is said that Y_n converges in probability to a number a , denoted $Y_n \xrightarrow{p} a$, if for any $\varepsilon > 0$ it follows that

$$\lim_{n \rightarrow \infty} \text{Prob}\{|Y_n - a| \geq \varepsilon\} = 0.$$

Convergence in probability can be also considered for a sequence $\mathbf{Y}_n \in \mathbb{R}^m$, $n = 1, \dots$, of random vectors. That is, \mathbf{Y}_n converges in probability to \mathbf{a} if for any $\varepsilon > 0$,

$$\lim_{n \rightarrow \infty} \text{Prob}\{\|\mathbf{Y}_n - \mathbf{a}\| \geq \varepsilon\} = 0.$$

It is straightforward to show that \mathbf{Y}_n converges in probability to \mathbf{a} iff its every component Y_{in} converges in probability to a_i , $i = 1, \dots, m$.

Law of Large Numbers (LLN) can be proved by using Chebishev inequality. Let X be a nonnegative valued random variable. Then for any $\varepsilon > 0$ we have

$$\text{Prob}(X \geq \varepsilon) = \mathbb{E}[\mathbf{1}_{[\varepsilon, \infty)}(X)] \leq \mathbb{E}[\varepsilon^{-1}X] = \varepsilon^{-1}\mathbb{E}[X],$$

where $\mathbf{1}_{[\varepsilon, \infty)}(x) = 0$ if $x < \varepsilon$ and $\mathbf{1}_{[\varepsilon, \infty)}(x) = 1$ if $x \geq \varepsilon$. The above inequality sometimes is called Markov inequality. Now let X be a random variable with finite second order moment, i.e., $\mathbb{E}[X^2] < \infty$. By taking $Y = (X - \mu)^2$, where $\mu = \mathbb{E}[X]$, we obtain from Markov inequality the following Chebishev inequality:

$$\text{Prob}\{|X - \mu| \geq \varepsilon\} = \text{Prob}\{(X - \mu)^2 \geq \varepsilon^2\} \leq \varepsilon^{-2}\mathbb{E}[(X - \mu)^2] = \varepsilon^{-2}\text{Var}(X).$$

It follows that if Y_n is a sequence of random variables such that $\mathbb{E}[Y_n] = \mu$, for all n , and $\text{Var}(Y_n)$ tends to zero as $n \rightarrow \infty$, then then for any $\varepsilon > 0$,

$$\text{Prob}\{|Y_n - \mu| \geq \varepsilon\} \leq \varepsilon^{-2} \text{Var}(Y_n) \rightarrow 0.$$

This implies that $Y_n \xrightarrow{p} \mu$. In particular, if X_1, \dots, X_n is iid with $\mu = \mathbb{E}[X_i]$ and $\sigma^2 = \text{Var}(X_i)$, then $\text{Var}(\bar{X}) = \sigma^2/n \rightarrow 0$, and hence $\bar{X} \xrightarrow{p} \mu$ as $n \rightarrow \infty$.

The convergence of \bar{X} to μ in probability is referred to as the (weak) Law of Large Numbers (WLLN). The stronger version of LLN is that \bar{X} converges to μ with probability one (w.p.1), provided the mean μ is well defined and finite. Note that convergence w.p.1 implies convergence in probability.

In Calculus the notation $y_n = o(x_n)$ is used to denote that if x_n and y_n are sequences of (deterministic) numbers, then y_n/x_n tends to zero as $n \rightarrow \infty$. The notation $y_n = O(x_n)$ means that there is a constant $C \geq 0$ such that $|y_n| \leq C|x_n|$ for all n . Now let X_n and Y_n be two sequences of random numbers. For random numbers counterparts of $o(\cdot)$ and $O(\cdot)$ are defined as follows. The notation $Y_n = o_p(X_n)$ means that $Y_n/X_n \xrightarrow{p} 0$ as $n \rightarrow \infty$. Usually it is used when X_n is deterministic. In particular $Y_n = o_p(1)$ means that $Y_n \xrightarrow{p} 0$. It is said that Y_n is *bounded in probability* if for any $\varepsilon > 0$ there exists $c > 0$ such that $\text{Prob}\{|Y_n| > c\} \leq \varepsilon$ for all n . The notation $Y_n = O_p(X_n)$ means that Y_n/X_n is bounded in probability. These notations $o_p(\cdot)$ and $O_p(\cdot)$ can be viewed as probabilistic analogues of their deterministic counterparts $o(\cdot)$ and $O(\cdot)$ and have similar properties. For example if $X_n = o_p(1)$ and $Y_n = O_p(1)$, then $X_n Y_n = o_p(1)$.

Recall that X_n converges in distribution to a random variable X , denoted $X_n \rightsquigarrow X$, if for any number x such that $\text{Prob}\{X = x\} = 0$ it follows that

$$\lim_{n \rightarrow \infty} \text{Prob}\{X_n \leq x\} = \text{Prob}\{X \leq x\}. \quad (6.1)$$

Note that condition $\text{Prob}\{X = x\} = 0$ means that the cumulative distribution function (cdf) $F(x) = \text{Prob}(X \leq x)$ of X is continuous at x , and condition (6.1) means that $\lim_{n \rightarrow \infty} F_n(x) = F(x)$, where $F_n(\cdot)$ is the cdf of X_n . That is, X_n converges in distribution to X if the cdf of X_n converges to the cdf of X at every point where the cdf of X is continuous.

A sequence of random vectors $\mathbf{X}_n \in \mathbb{R}^m$ converges in distribution to a random vector \mathbf{X} if

$$\lim_{n \rightarrow \infty} \text{Prob}\{\mathbf{X}_n \in A\} = \text{Prob}\{\mathbf{X} \in A\}$$

for any rectangular set $A = \{\mathbf{x} : a_i \leq x_i \leq b_i, i = 1, \dots, m\}$ such that probability of \mathbf{X} to be on the boundary of A is zero.

Proposition 6.1 *If $X_n \rightsquigarrow X$, then $X_n = O_p(1)$, i.e., if X_n converges in distribution, then X_n is bounded in probability.*

Proof. Let $F_n(x) = \text{Prob}(X_n \leq x)$ and $F(x) = \text{Prob}(X \leq x)$ be cumulative distribution functions of X_n and X , respectively, and $\varepsilon > 0$. Recall that X_n converges in distribution to X iff $\lim_{n \rightarrow \infty} F_n(x) = F(x)$ for every $x \in \mathbb{R}$ such that $F(\cdot)$ is continuous at x . Therefore we have that $F_n(x) \rightarrow F(x)$ provided that F is continuous at x . Since $F(x) \rightarrow 1$ as $x \rightarrow +\infty$ and $F(x) \rightarrow 0$ as $x \rightarrow -\infty$, there exists a constant c_0 such that $F(c_0) > 1 - \varepsilon$ and $F(-c_0) < \varepsilon$. Moreover, since a monotonically nondecreasing function can have only a countable number of discontinuous points, we can choose this constant c_0 such that F is continuous at c_0 . It follows that there exists N such that $F_n(c_0) \geq 1 - 2\varepsilon$ and $F_n(c_0) \leq 2\varepsilon$ for all $n \geq N$. That is, $\text{Prob}(|X_n| \geq c_0) \leq 4\varepsilon$. Now for every k there is a constant c_k such that $\text{Prob}(|X_k| \geq c_k) \leq \varepsilon$. Then for $c = \max\{c_0, c_1, \dots, c_N\}$ we have that $\text{Prob}(|X_n| \geq c) \leq 4\varepsilon$ for all $n \in \mathbb{N}$. That is, for any $\varepsilon > 0$ there is c such that $\text{Prob}(|X_n| \geq c) \leq 4\varepsilon$ for all n . This shows that X_n is bounded in probability. \square

Theorem 6.1 (Slutsky's theorem) *If $X_n \rightsquigarrow X$ and $Y_n \xrightarrow{P} 0$, then $X_n + Y_n \rightsquigarrow X$.*

Proof. Consider the cdf $F(x) = \text{Prob}\{X \leq x\}$ of X . Let x be such that $F(\cdot)$ is continuous at x . We need to show that $\text{Prob}\{X_n + Y_n \leq x\}$ tends to $F(x)$ as $n \rightarrow \infty$. For $\varepsilon > 0$ we have

$$\begin{aligned} \text{Prob}(X_n + Y_n \leq x) &= \text{Prob}(X_n + Y_n \leq x, |Y_n| \leq \varepsilon) + \text{Prob}(X_n + Y_n \leq x, |Y_n| > \varepsilon) \\ &\leq \text{Prob}(X_n \leq x + \varepsilon) + \text{Prob}(|Y_n| > \varepsilon). \end{aligned}$$

Since $Y_n \xrightarrow{P} 0$ we have that $\text{Prob}(|Y_n| > \varepsilon)$ tends to zero. Moreover let $\varepsilon > 0$ be such that $F(\cdot)$ is continuous at $x + \varepsilon$. Then since $X_n \rightsquigarrow X$, we have that $\text{Prob}(X_n \leq x + \varepsilon)$ tends to $F(x + \varepsilon)$. It follows that

$$\limsup_{n \rightarrow \infty} \text{Prob}(X_n + Y_n \leq x) \leq F(x + \varepsilon).$$

Note that since $F(x)$ is monotonically nondecreasing, the set of points where it is discontinuous is countable. Therefore we can choose a sequence $\varepsilon_n \downarrow 0$ such that $F(\cdot)$ is continuous at $x + \varepsilon_n$. By continuity of $F(\cdot)$ at x it follows that

$$\limsup_{n \rightarrow \infty} \text{Prob}(X_n + Y_n \leq x) \leq F(x).$$

In a similar way it is possible to show that

$$\liminf_{n \rightarrow \infty} \text{Prob}(X_n + Y_n \leq x) \geq F(x).$$

It follows that $\text{Prob}\{X_n + Y_n \leq x\}$ tends to $F(x)$. □

The concept of ‘bounded in probability’ can be extended to a sequence $\mathbf{Y}_n \in \mathbb{R}^m$ of random vectors. That is, \mathbf{Y}_n is bounded in probability if for any $\varepsilon > 0$ there is a bounded set $A \subset \mathbb{R}^m$ such that $\text{Prob}\{\mathbf{Y}_n \notin A\} \leq \varepsilon$ for all n . It is not difficult to show that \mathbf{Y}_n is bounded in probability iff its every component Y_{in} is bounded in probability.

Slutsky's theorem also can be extended to random vectors. That is, if \mathbf{X}_n converges in distribution to \mathbf{X} and \mathbf{Y}_n converges in probability to $\mathbf{0}$, then $\mathbf{X}_n + \mathbf{Y}_n$ converges in distribution to \mathbf{X} .

Theorem 6.2 (Delta theorem) *Let \mathbf{X}_n be a sequence of $m \times 1$ random vectors and $\mathbf{g} : \mathbb{R}^m \rightarrow \mathbb{R}^k$ be a function. Suppose that $\lambda_n(\mathbf{X}_n - \boldsymbol{\mu}) \rightsquigarrow \mathbf{Z}$, where $\boldsymbol{\mu} \in \mathbb{R}^m$ and $\lambda_n \rightarrow \infty$, and that $\mathbf{g}(\cdot)$ is differentiable at $\boldsymbol{\mu}$ with $\nabla \mathbf{g}(\boldsymbol{\mu}) = \partial \mathbf{g}(\boldsymbol{\mu}) / \partial \mathbf{x}$ being the $m \times k$ matrix of partial derivatives (Jacobian matrix). Then*

$$\lambda_n(\mathbf{g}(\mathbf{X}_n) - \mathbf{g}(\boldsymbol{\mu})) \rightsquigarrow [\nabla \mathbf{g}(\boldsymbol{\mu})]' \mathbf{Z}. \quad (6.2)$$

Proof. Since $\mathbf{g}(\cdot)$ is differentiable at $\boldsymbol{\mu}$ we have that

$$\mathbf{g}(\mathbf{x}) - \mathbf{g}(\boldsymbol{\mu}) = [\nabla \mathbf{g}(\boldsymbol{\mu})]'(\mathbf{x} - \boldsymbol{\mu}) + \mathbf{r}(\mathbf{x}),$$

where $\boldsymbol{\varepsilon}(\mathbf{x}) = \mathbf{r}(\mathbf{x}) / \|\mathbf{x} - \boldsymbol{\mu}\|$ tends to $\mathbf{0}$ as $\mathbf{x} \rightarrow \boldsymbol{\mu}$. Hence

$$\lambda_n(\mathbf{g}(\mathbf{X}_n) - \mathbf{g}(\boldsymbol{\mu})) = [\nabla \mathbf{g}(\boldsymbol{\mu})]'[\lambda_n(\mathbf{X}_n - \boldsymbol{\mu})] + \boldsymbol{\varepsilon}(\mathbf{X}_n)[\lambda_n \|\mathbf{X}_n - \boldsymbol{\mu}\|]. \quad (6.3)$$

Now since $\lambda_n(\mathbf{X}_n - \boldsymbol{\mu})$ converges in distribution, it follows that $\lambda_n(\mathbf{X}_n - \boldsymbol{\mu})$ is bounded in probability. Moreover since $\lambda_n \rightarrow \infty$ it follows that $\mathbf{X}_n \xrightarrow{P} \boldsymbol{\mu}$. Hence $\boldsymbol{\varepsilon}(\mathbf{X}_n) \xrightarrow{P} \mathbf{0}$, and thus

$\varepsilon(\mathbf{X}_n)[\lambda_n \|\mathbf{X}_n - \boldsymbol{\mu}\|] \xrightarrow{P} \mathbf{0}$. By Slutsky's theorem the convergence (6.2) follows from (6.3). \square

In particular it follows that if in addition to the assumptions of Theorem 6.2, $\sqrt{n}(\mathbf{X}_n - \boldsymbol{\mu})$ converges in distribution to normal $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$, then $\sqrt{n}(\mathbf{g}(\mathbf{X}_n) - \mathbf{g}(\boldsymbol{\mu}))$ converges in distribution to normal with zero mean and covariance matrix $[\nabla \mathbf{g}(\boldsymbol{\mu})]' \boldsymbol{\Sigma} [\nabla \mathbf{g}(\boldsymbol{\mu})]$. For $k = 1$, i.e., when $g : \mathbb{R}^m \rightarrow \mathbb{R}$ is a real valued function, $\nabla g(\boldsymbol{\mu})$ becomes the gradient $\nabla g(\boldsymbol{\mu}) = (\partial g(\boldsymbol{\mu})/\partial x_1, \dots, \partial g(\boldsymbol{\mu})/\partial x_m)'$ and $[\nabla g(\boldsymbol{\mu})]' \boldsymbol{\Sigma} [\nabla g(\boldsymbol{\mu})]$ becomes the asymptotic variance of $\sqrt{n}(g(\mathbf{X}_n) - g(\boldsymbol{\mu}))$.

Example 6.1 Let X_n and Y_n be two independent sequences of random variables such that $\sqrt{n}(X_n - \mu_x) \rightsquigarrow \mathcal{N}(0, \sigma_x^2)$ and $\sqrt{n}(Y_n - \mu_y) \rightsquigarrow \mathcal{N}(0, \sigma_y^2)$, $\mu_y \neq 0$. Let us find the asymptotic distribution of (V_n, W_n) , where $V_n = X_n Y_n$ and $W_n = X_n/Y_n$. Consider $\mathbf{g}(x, y) = (xy, x/y)$. Note that $\mathbf{g}(X_n, Y_n) = (V_n, W_n)$. By Delta Theorem we have that $\sqrt{n} \begin{bmatrix} V_n - \mu_x \mu_y \\ W_n - \mu_x / \mu_y \end{bmatrix}$ converges in distribution to normal $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ with

$$\boldsymbol{\Sigma} = \begin{bmatrix} \mu_y & \mu_x \\ 1/\mu_y & -\mu_x/\mu_y^2 \end{bmatrix} \begin{bmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{bmatrix} \begin{bmatrix} \mu_y & 1/\mu_y \\ \mu_x & -\mu_x/\mu_y^2 \end{bmatrix}.$$

That is, elements of the asymptotic covariance matrix $\boldsymbol{\Sigma}$ are: $\sigma_{11} = \mu_y^2 \sigma_x^2 + \mu_x^2 \sigma_y^2$, $\sigma_{22} = \sigma_x^2 / \mu_y^2 + (\mu_x^2 / \mu_y^4) \sigma_y^2$, $\sigma_{12} = \sigma_x^2 - (\mu_x / \mu_y)^2 \sigma_y^2$. In particular, if $\mu_x = \mu_y$ and $\sigma_x^2 = \sigma_y^2$, then $\sigma_{12} = 0$. In that case $V_n = X_n Y_n$ and $W_n = X_n / Y_n$ are asymptotically independent. \square

Delta method can be extended to high order terms. For example suppose that $g : \mathbb{R}^m \rightarrow \mathbb{R}$ is twice continuously differentiable and $\nabla g(\boldsymbol{\mu}) = \mathbf{0}$. Then the right hand side of (6.2) degenerates to 0. Let \mathbf{H} be the $m \times m$ Hessian matrix of second order partial derivatives at $\mathbf{x} = \boldsymbol{\mu}$, i.e., $H_{ij} = \frac{\partial^2 g(\boldsymbol{\mu})}{\partial x_i \partial x_j}$, $i, j = 1, \dots, m$. The second order expansion of $g(\cdot)$ at $\boldsymbol{\mu}$ is

$$g(\mathbf{x}) - g(\boldsymbol{\mu}) = \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \mathbf{H}(\mathbf{x} - \boldsymbol{\mu}) + r(\mathbf{x}),$$

where the remainder $r(\mathbf{x})$ is of order $o(\|\mathbf{x} - \boldsymbol{\mu}\|^2)$, i.e., $r(\mathbf{x})/\|\mathbf{x} - \boldsymbol{\mu}\|^2$ tends to 0 as $\mathbf{x} \rightarrow \boldsymbol{\mu}$. Suppose further that $\sqrt{n}(\mathbf{X}_n - \boldsymbol{\mu}) \rightsquigarrow \mathbf{Z}$. Then

$$n(g(\mathbf{X}_n) - g(\boldsymbol{\mu})) \rightsquigarrow \frac{1}{2} \mathbf{Z}' \mathbf{H} \mathbf{Z}. \quad (6.4)$$

That is, $2n(g(\mathbf{X}_n) - g(\boldsymbol{\mu}))$ converges in distribution to the quadratic form $Q = \mathbf{Z}' \mathbf{H} \mathbf{Z}$.

7 Exponential family of distributions

It is said that \mathbf{X} is distributed according to the *exponential family* (in the canonical form) if its probability density function (pdf) is of the form

$$f(\mathbf{x}, \boldsymbol{\theta}) = \exp \left[\sum_{i=1}^k \theta_i T_i(\mathbf{x}) - A(\boldsymbol{\theta}) \right] h(\mathbf{x}), \quad (7.1)$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)' \in \Theta$ is vector of parameters with

$$\Theta = \left\{ \boldsymbol{\theta} : \int \exp \left[\sum_{i=1}^k \theta_i T_i(\mathbf{x}) \right] h(\mathbf{x}) d\mathbf{x} < \infty \right\}.$$

Let us show that for $T_j = T_j(\mathbf{X})$ and $\mathbb{E}_{\boldsymbol{\theta}}(T_j) = \int T_j(\mathbf{x}) f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x}$,

$$\mathbb{E}_{\boldsymbol{\theta}}(T_j) = \frac{\partial}{\partial \theta_j} A(\boldsymbol{\theta}), \quad (7.2)$$

$$\text{Cov}(T_j, T_\ell) = \frac{\partial^2}{\partial\theta_j\partial\theta_\ell} A(\boldsymbol{\theta}). \quad (7.3)$$

Indeed, we have that $\int f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} = 1$ for all $\boldsymbol{\theta} \in \Theta$. Let $\boldsymbol{\theta}$ be an interior point of Θ , and hence the expectation and differentiation can be interchanged. We have that $\frac{\partial}{\partial\theta_j} \int f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} = 0$ and

$$\frac{\partial}{\partial\theta_j} f(\mathbf{x}, \boldsymbol{\theta}) = \left[T_j(\mathbf{x}) - \frac{\partial}{\partial\theta_j} A(\boldsymbol{\theta}) \right] f(\mathbf{x}, \boldsymbol{\theta}),$$

and hence

$$\begin{aligned} 0 &= \frac{\partial}{\partial\theta_j} \int f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} = \int \frac{\partial}{\partial\theta_j} f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} \\ &= \mathbb{E}_\theta \left[T_j - \frac{\partial}{\partial\theta_j} A(\boldsymbol{\theta}) \right] = \mathbb{E}_\theta(T_j) - \frac{\partial}{\partial\theta_j} A(\boldsymbol{\theta}). \end{aligned}$$

It follows that $\mathbb{E}_\theta(T_j) = \frac{\partial}{\partial\theta_j} A(\boldsymbol{\theta})$. The other equation follows in a similar way from $\frac{\partial^2}{\partial\theta_j\partial\theta_\ell} \int f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} = 0$.

8 Point estimation

8.1 Maximum likelihood method

Consider a parametric family of distributions defined by probability density functions (pdf) $f(\mathbf{x}, \boldsymbol{\theta})$, $\mathbf{x} \in \mathbb{R}^m$, with parameter vector $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^k$. Given an iid sample $\mathbf{X}_1, \dots, \mathbf{X}_N$, the Maximum Likelihood (ML) estimator of $\boldsymbol{\theta}$ is the maximizer $\hat{\boldsymbol{\theta}}_n$ of the likelihood function $L_N(\boldsymbol{\theta}) = \prod_{i=1}^N f(\mathbf{X}_i, \boldsymbol{\theta})$ over $\boldsymbol{\theta} \in \Theta$. Note that both $L_N(\boldsymbol{\theta})$ and $\hat{\boldsymbol{\theta}}_N$ are functions of the sample, this is suppressed in the notation. Since $\log x$ is monotonically increasing for $x > 0$, this can be written as

$$\hat{\boldsymbol{\theta}}_N \in \arg \max_{\boldsymbol{\theta} \in \Theta} \log L_N(\boldsymbol{\theta}). \quad (8.1)$$

Note that such maximizer may not exist or could be not unique. We assume that the random sample is an iid replication of random vector \mathbf{X} having pdf $g(\mathbf{x})$, written $\mathbf{X} \sim g(\cdot)$, i.e., each \mathbf{X}_i has pdf $g(\cdot)$. In particular if $g(\cdot) = f(\cdot, \boldsymbol{\theta}^*)$ for some $\boldsymbol{\theta}^* \in \Theta$, we say that the model is *correctly specified*. It is said that the model is *identified* at $\boldsymbol{\theta}^*$ if $f(\cdot, \boldsymbol{\theta}) = f(\cdot, \boldsymbol{\theta}^*)$, $\boldsymbol{\theta} \in \Theta$, implies that $\boldsymbol{\theta} = \boldsymbol{\theta}^*$. That is, $\boldsymbol{\theta}^*$ is the unique value of the parameter vector which defines the model.

Since $\log L_N(\boldsymbol{\theta}) = \sum_{i=1}^N \log f(\mathbf{X}_i, \boldsymbol{\theta})$, it follows by the LLN that for a given $\boldsymbol{\theta}$ the average $N^{-1} \log L_N(\boldsymbol{\theta})$ converges w.p.1 as $N \rightarrow \infty$ to

$$\mathbb{E}_g[\log f(\mathbf{X}, \boldsymbol{\theta})] = \int [\log f(\mathbf{x}, \boldsymbol{\theta})] g(\mathbf{x}) d\mathbf{x},$$

provided this expectation is well defined and finite. The notation \mathbb{E}_g emphasizes that the expectation is taken with respect to the distribution of the sample defined by the pdf $g(\cdot)$. It is natural then to expect that the ML estimator $\hat{\boldsymbol{\theta}}_N$ will converge w.p.1 to a maximizer of $\mathbb{E}_g[\log f(\mathbf{X}, \boldsymbol{\theta})]$ over $\boldsymbol{\theta} \in \Theta$. And indeed it is possible to prove that such converges holds under certain regularity conditions. In order to understand what such maximizer is, we need the following inequality.

Theorem 8.1 (Jensen inequality) *Let $\phi: \mathbb{R}^m \rightarrow \mathbb{R}$ be a convex function and \mathbf{X} be an $m \times 1$ random vector having mean $\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}]$. Then*

$$\mathbb{E}[\phi(\mathbf{X})] \geq \phi(\boldsymbol{\mu}). \quad (8.2)$$

Proof. Since $\phi(\cdot)$ is convex we have that there exists $\gamma \in \mathbb{R}^m$ such that

$$\phi(\mathbf{x}) \geq \phi(\boldsymbol{\mu}) + \gamma'(\mathbf{x} - \boldsymbol{\mu})$$

for any $\mathbf{x} \in \mathbb{R}^m$ (vector γ is called subgradient of ϕ at $\boldsymbol{\mu}$). It follows that

$$\mathbb{E}[\phi(\mathbf{X})] \geq \phi(\boldsymbol{\mu}) + \mathbb{E}[\gamma'(\mathbf{X} - \boldsymbol{\mu})].$$

Since $\mathbb{E}[\gamma'(\mathbf{X} - \boldsymbol{\mu})] = \gamma'(\mathbb{E}[\mathbf{X}] - \boldsymbol{\mu}) = 0$, the inequality (8.2) follows. \square

Kullback-Leibler divergence of pdf $f(\cdot)$ from pdf $g(\cdot)$ is defined as

$$D(g\|f) := \int \left[\log \frac{g(\mathbf{x})}{f(\mathbf{x})} \right] g(\mathbf{x}) d\mathbf{x} = \mathbb{E}_g \left[\log \frac{g(\mathbf{X})}{f(\mathbf{X})} \right] = -\mathbb{E}_g \left[\log \frac{f(\mathbf{X})}{g(\mathbf{X})} \right].$$

Since $-\log x$ is a convex function we have by Jensen inequality

$$\begin{aligned} D(g\|f) &= -\mathbb{E}_g \left[\log \frac{f(\mathbf{X})}{g(\mathbf{X})} \right] \geq -\log \mathbb{E}_g \left[\frac{f(\mathbf{X})}{g(\mathbf{X})} \right] \\ &= -\log \int \frac{f(\mathbf{x})}{g(\mathbf{x})} g(\mathbf{x}) d\mathbf{x} = -\log \int f(\mathbf{x}) d\mathbf{x} = -\log 1 = 0. \end{aligned}$$

That is, $D(g\|f) \geq 0$ and $D(g\|f) = 0$ iff $f = g$.

Since

$$D(g(\cdot)\|f(\cdot, \boldsymbol{\theta})) = \mathbb{E}_g [\log g(\mathbf{X})] - \mathbb{E}_g [\log f(\mathbf{X}, \boldsymbol{\theta})],$$

we have that maximizing $\mathbb{E}_g [\log f(\mathbf{X}, \boldsymbol{\theta})]$, over $\boldsymbol{\theta} \in \Theta$, is equivalent to minimizing the Kullback-Leibler divergence of $f(\cdot, \boldsymbol{\theta})$ from $g(\cdot)$. In particular, if the model is correctly specified, i.e., $g(\cdot) = f(\cdot, \boldsymbol{\theta}^*)$ for some $\boldsymbol{\theta}^* \in \Theta$, then $\boldsymbol{\theta}^*$ is a maximizer of $\mathbb{E}_{\theta^*} [\log f(\mathbf{X}, \boldsymbol{\theta})]$, over $\boldsymbol{\theta} \in \Theta$, where the notation \mathbb{E}_{θ^*} emphasizes that the expectation is taken with respect to the distribution $g(\cdot) = f(\cdot, \boldsymbol{\theta}^*)$. That is

$$\boldsymbol{\theta}^* \in \arg \max_{\boldsymbol{\theta} \in \Theta} \left\{ \mathbb{E}_{\theta^*} [\log f(\mathbf{X}, \boldsymbol{\theta})] = \int [\log f(\mathbf{x}, \boldsymbol{\theta})] f(\mathbf{x}, \boldsymbol{\theta}^*) d\mathbf{x} \right\}.$$

It follows that if the model is correctly specified and identified at $\boldsymbol{\theta}^*$ and some regularity conditions are satisfied, then the ML estimator $\hat{\boldsymbol{\theta}}_N$ converges w.p.1 to $\boldsymbol{\theta}^*$. In that case it is said that $\hat{\boldsymbol{\theta}}_N$ is a *consistent* estimator of $\boldsymbol{\theta}^*$.

8.1.1 Asymptotic distribution of the ML estimators

Let $\mathbf{X} \sim f(\mathbf{x}, \boldsymbol{\theta})$, $\boldsymbol{\theta} \in \mathbb{R}^k$, be a random vector. The following $k \times k$ matrix is called (Fisher) information matrix

$$\mathbf{I}(\boldsymbol{\theta}) := \mathbb{E}_{\boldsymbol{\theta}} \left\{ \left[\frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{X}, \boldsymbol{\theta}) \right] \left[\frac{\partial}{\partial \boldsymbol{\theta}} \log f(\mathbf{X}, \boldsymbol{\theta}) \right]' \right\}. \quad (8.3)$$

The notation $\mathbb{E}_{\boldsymbol{\theta}}$ emphasises that the expectation is taken with respect to the distribution $f(\cdot, \boldsymbol{\theta})$ of \mathbf{X} . Note that $\mathbf{I}(\boldsymbol{\theta})$ is a function of $\boldsymbol{\theta}$.

Let us show that

$$\mathbf{I}(\boldsymbol{\theta}) = -\mathbb{E}_{\boldsymbol{\theta}} \left\{ \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log f(\mathbf{X}, \boldsymbol{\theta}) \right\}. \quad (8.4)$$

We need to show that

$$\mathbb{E}_\theta \left\{ \frac{\partial \log f(\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_i} \frac{\partial \log f(\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_j} \right\} = -\mathbb{E}_\theta \left\{ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(\mathbf{X}, \boldsymbol{\theta}) \right\}, \quad (8.5)$$

$i, j = 1, \dots, k$. We have that

$$\mathbb{E}_\theta \left\{ \frac{\partial \log f(\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_i} \right\} = \mathbb{E}_\theta \left\{ \frac{\partial f(\mathbf{X}, \boldsymbol{\theta}) / \partial \theta_i}{f(\mathbf{X}, \boldsymbol{\theta})} \right\} = \int \frac{\partial f(\mathbf{x}, \boldsymbol{\theta}) / \partial \theta_i}{f(\mathbf{x}, \boldsymbol{\theta})} f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} = \int \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i} d\mathbf{x}.$$

Suppose now that

$$\int \frac{\partial f(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i} d\mathbf{x} = \frac{\partial}{\partial \theta_i} \int f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x}, \quad (8.6)$$

i.e., the operations of differentiation with respect to θ_i and integration with respect to \mathbf{x} can be interchanged. Then the right hand side of (8.6) is 0, since $\int f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} = 1$ for all $\boldsymbol{\theta}$. It follows that

$$\mathbb{E}_\theta \left\{ \frac{\partial \log f(\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_i} \right\} = \int \frac{\partial \log f(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i} f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} = 0, \quad (8.7)$$

for all $\boldsymbol{\theta}$ and hence

$$\frac{\partial}{\partial \theta_j} \int \frac{\partial \log f(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i} f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} = 0. \quad (8.8)$$

By taking the derivative, in the left hand side of (8.8), inside the integral we obtain

$$\begin{aligned} 0 &= \int \frac{\partial}{\partial \theta_j} \left[\frac{\partial \log f(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i} f(\mathbf{x}, \boldsymbol{\theta}) \right] d\mathbf{x} \\ &= \int \frac{\partial^2 \log f(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} + \int \frac{\partial \log f(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i} \frac{\partial \log f(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_j} f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} \\ &= \mathbb{E}_\theta \left\{ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(\mathbf{X}, \boldsymbol{\theta}) \right\} + \mathbb{E}_\theta \left\{ \frac{\partial \log f(\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_i} \frac{\partial \log f(\mathbf{X}, \boldsymbol{\theta})}{\partial \theta_j} \right\}, \end{aligned} \quad (8.9)$$

and hence (8.5) follows. \square

Remark 8.1 The above derivations are based on the interchangeability property that the operations of differentiation with respect to θ_i and integration with respect to \mathbf{x} can be interchanged. We used it twice, in (8.6) and again in (8.9). As it is discussed below the interchangeability property (8.6) holds if $f(\mathbf{x}, \cdot)$ is differentiable and there is nonnegative valued function $K(\mathbf{x})$ such that $\mathbb{E}[K(\mathbf{X})] < \infty$ and

$$|f(\mathbf{x}, \boldsymbol{\theta}_1) - f(\mathbf{x}, \boldsymbol{\theta}_2)| \leq K(\mathbf{x}) \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|, \quad \boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^k,$$

i.e., if $f(\mathbf{x}, \cdot)$ is Lipschitz continuous with integrable Lipschitz constant. Similar condition is needed for $\partial f(\mathbf{x}, \boldsymbol{\theta}) / \partial \theta_i$, $i = 1, \dots, k$, in order to justify (8.9).

Let us discuss conditions ensuring that the expectation and differentiation can be interchanged. Let $g(x, \theta)$ be a real valued function of $x, \theta \in \mathbb{R}$. Suppose that $g(x, \theta)$ is differentiable in θ . We would like to verify that

$$\frac{\partial}{\partial \theta} \mathbb{E}[g(X, \theta)] = \mathbb{E} \left[\frac{\partial}{\partial \theta} g(X, \theta) \right],$$

where the expectation is with respect to distribution of random variable X . We have

$$\frac{\partial}{\partial \theta} \mathbb{E}[g(X, \theta)] = \lim_{h \rightarrow 0} \frac{\mathbb{E}[g(X, \theta + h)] - \mathbb{E}[g(X, \theta)]}{h} = \lim_{h \rightarrow 0} \mathbb{E} \left[\frac{g(X, \theta + h) - g(X, \theta)}{h} \right].$$

In order to interchange the limit and the expectation (integration) we can use the Lebesgue Dominated Convergence Theorem: if $f_n, g : \Omega \rightarrow \mathbb{R}$ are such that $|f_n| \leq g$, $\int_{\Omega} g dP < \infty$ and $f_n(\omega) \rightarrow f(\omega)$ for a.e. $\omega \in \Omega$, then $\int_{\Omega} f_n dP \rightarrow \int_{\Omega} f dP$.

That is, suppose that there is function $K(x) \geq 0$ such that $\mathbb{E}[K(X)] < \infty$ and for all h ,

$$|g(X, \theta + h) - g(X, \theta)| \leq K(X)|h|.$$

Then by the Lebesgue Dominated Convergence Theorem, the limit and the expectation (integration) can be interchanged and hence

$$\frac{\partial}{\partial \theta} \mathbb{E}[g(X, \theta)] = \mathbb{E} \left[\lim_{h \rightarrow 0} \frac{g(X, \theta + h) - g(X, \theta)}{h} \right] = \mathbb{E} \left[\frac{\partial}{\partial \theta} g(X, \theta) \right].$$

□

Let us show that the information matrix $\mathbf{I}(\theta)$ is positive semidefinite. We have that, for $\mathbf{a} \in \mathbb{R}^k$,

$$\mathbf{a}' \mathbf{I}(\theta) \mathbf{a} = \sum_{i,j=1}^k a_i a_j I_{ij}(\theta),$$

where

$$I_{ij}(\theta) = \mathbb{E}_{\theta} \left\{ \frac{\partial \log f(\mathbf{X}, \theta)}{\partial \theta_i} \frac{\partial \log f(\mathbf{X}, \theta)}{\partial \theta_j} \right\},$$

and hence

$$a_i a_j I_{ij}(\theta) = \mathbb{E}_{\theta} \left\{ \left(a_i \frac{\partial \log f(\mathbf{X}, \theta)}{\partial \theta_i} \right) \left(a_j \frac{\partial \log f(\mathbf{X}, \theta)}{\partial \theta_j} \right) \right\}.$$

It follows that

$$\mathbf{a}' \mathbf{I}(\theta) \mathbf{a} = \mathbb{E}_{\theta} \left\{ \left[\sum_{i=1}^k a_i \frac{\partial \log f(\mathbf{X}, \theta)}{\partial \theta_i} \right]^2 \right\},$$

and hence $\mathbf{a}' \mathbf{I}(\theta) \mathbf{a} \geq 0$. □

Consider now the ML estimation procedure. Suppose that the model is correctly specified and let $\hat{\theta}_N$ be the ML estimator of the true parameter value θ^* . Assume that $\hat{\theta}_N$ is a consistent estimator of θ^* , i.e., $\hat{\theta}_N$ converges w.p.1 to θ^* . Suppose further that θ^* is an interior point of the set Θ . Since $\hat{\theta}_N$ is a consistent estimator of θ^* , it follows that $\hat{\theta}_N$ is in the interior of the set Θ for all N large enough. Then since $\hat{\theta}_N$ is a maximizer of $\log L_N(\theta)$, the following optimality condition holds

$$\frac{\partial}{\partial \theta} \left[\sum_{i=1}^N \log f(\mathbf{X}_i, \hat{\theta}_N) \right] = \mathbf{0}. \quad (8.10)$$

By the Mean Value Theorem we can write

$$\frac{\partial}{\partial \theta} \left[\sum_{i=1}^N \log f(\mathbf{X}_i, \hat{\theta}_N) \right] = \frac{\partial}{\partial \theta} \left[\sum_{i=1}^N \log f(\mathbf{X}_i, \theta^*) \right] + \left[\frac{\partial^2}{\partial \theta \partial \theta'} \sum_{i=1}^N \log f(\mathbf{X}_i, \tilde{\theta}_N) \right] (\hat{\theta}_N - \theta^*),$$

for some $\tilde{\theta}_N$ between $\hat{\theta}_N$ and θ^* . It follows that

$$\begin{aligned}\sqrt{N}(\hat{\theta}_N - \theta^*) &= -\sqrt{N} \left[\frac{\partial^2}{\partial \theta \partial \theta'} \sum_{i=1}^N \log f(\mathbf{X}_i, \tilde{\theta}_N) \right]^{-1} \left[\frac{\partial}{\partial \theta} \sum_{i=1}^N \log f(\mathbf{X}_i, \theta^*) \right] \\ &= - \left[\frac{1}{N} \frac{\partial^2}{\partial \theta \partial \theta'} \sum_{i=1}^N \log f(\mathbf{X}_i, \tilde{\theta}_N) \right]^{-1} \left[\frac{1}{\sqrt{N}} \frac{\partial}{\partial \theta} \sum_{i=1}^N \log f(\mathbf{X}_i, \theta^*) \right].\end{aligned}\quad (8.11)$$

Since $\hat{\theta}_N$, and hence $\tilde{\theta}_N$, converge to θ^* w.p.1, we have by the LLN that

$$\frac{1}{N} \frac{\partial^2}{\partial \theta \partial \theta'} \sum_{i=1}^N \log f(\mathbf{X}_i, \tilde{\theta}_N) = \frac{1}{N} \sum_{i=1}^N \frac{\partial^2}{\partial \theta \partial \theta'} \log f(\mathbf{X}_i, \tilde{\theta}_N)$$

converges to $-\mathbf{I}(\theta^*)$. Now note that

$$\mathbb{E}_\theta \left[\frac{\partial}{\partial \theta_i} \log f(\mathbf{X}, \theta) \right] = \mathbb{E}_\theta \left[\frac{\frac{\partial}{\partial \theta_i} f(\mathbf{X}, \theta)}{f(\mathbf{X}, \theta)} \right] = \int \frac{\partial}{\partial \theta_i} f(\mathbf{x}, \theta) d\mathbf{x} = \frac{\partial}{\partial \theta_i} \int f(\mathbf{x}, \theta) d\mathbf{x} = 0.$$

Therefore by the CLT we have that $\frac{1}{\sqrt{N}} \frac{\partial}{\partial \theta} \sum_{i=1}^N \log f(\mathbf{X}_i, \theta^*)$ converges in distribution to normal with zero mean vector and covariance matrix $\mathbf{I}(\theta^*)$. Together with (8.11) this implies that

$$\sqrt{N}(\hat{\theta}_N - \theta^*) \rightsquigarrow \mathcal{N}(\mathbf{0}, \mathbf{I}(\theta^*)^{-1}). \quad (8.12)$$

That is, the ML estimator $\hat{\theta}_N$ has approximately normal distribution with mean θ^* and covariance matrix $N^{-1} \mathbf{I}(\theta^*)^{-1}$. \square

Remark 8.2 The above derivations of the basic result (8.12) involve several assumptions (regularity conditions). The asymptotic result (8.12) is local, it is based on the second order approximation of the likelihood function at the true value θ^* . So the MLE estimator should be consistent in order to justify such approximations. In order to apply necessary condition (8.10), the MLE should be an interior point of the set Θ , i.e., should not be on the boundary of the set Θ . If θ^* is a boundary point of Θ , then the asymptotics of the MLE is different. We also needed the interchangeability property, that the operations of integration with respect to \mathbf{x} and differentiation with respect to θ can be interchanged (see Remark 8.1). \square

Example 8.1 Suppose that X_1, \dots, X_N are iid having uniform distribution on the interval $[0, \theta]$, $\theta > 0$, with pdf $f(x, \theta) = 1/\theta$ for $x \in [0, \theta]$, and $f(x, \theta) = 0$ otherwise. Hence the likelihood function is $L_N(\theta) = 1/\theta^N$ for $X_{(N)} \leq \theta$, where $X_{(N)} = \max\{X_1, \dots, X_N\}$. Since $L_N(\theta)$ is monotonically decreasing with increase of θ , the MLE is given by the smallest possible value of θ which is $X_{(N)}$. That is $X_{(N)}$ is the MLE of θ .

The cdf of X_i is $F(x) = x/\theta$ for $x \in [0, \theta]$. Then the cdf of $N[\theta - X_{(N)}]$, for $x \in [0, N\theta]$, is

$$\begin{aligned}\text{Prob}(N[\theta - X_{(N)}] \leq x) &= \text{Prob}(X_{(N)} \geq \theta - x/N) = 1 - \text{Prob}(X_{(N)} < \theta - x/N) \\ &= 1 - \text{Prob}(X_i < \theta - x/N, i = 1, \dots, N) \\ &= 1 - \prod_{i=1}^N \text{Prob}(X_i < \theta - x/N) \\ &= 1 - [F(\theta - x/N)]^N = 1 - (1 - x/(N\theta))^N.\end{aligned}$$

Furthermore

$$\lim_{N \rightarrow \infty} (1 - x/(N\theta))^N = e^{-x/\theta}.$$

It follows that the cdf of $N[\theta - X_{(N)}]$ converges to $1 - e^{-x/\theta}$. This implies that $N[\theta - X_{(N)}]$ converges in distribution to exponential $\exp(\lambda)$ with $\lambda = 1/\theta$. Note that the situation here is not standard, the optimality equation (8.10) is not applicable here. Also the asymptotic variance is of order $O(N^{-2})$ rather than $O(N^{-1})$ as in the standard case. \square

8.2 Cramér - Rao lower bound

Let $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_N)$ be an iid sample from $f(\mathbf{x}, \theta)$, $\theta \in \mathbb{R}$, and $T(\mathbf{X})$ be a statistic, i.e., $T(\mathbf{X})$ is a function of \mathbf{X} . Note that $f(\mathbf{x}, \theta) = \prod_{j=1}^N f_j(\mathbf{x}_j, \theta)$ is the pdf of $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_N)$, where $f_j(\mathbf{x}_j, \theta)$ is pdf of \mathbf{X}_j . Since the sample is iid, pdfs $f_j(\cdot, \theta)$ are the same for all $j = 1, \dots, N$.

Then under some regularity conditions

$$\text{Var}_\theta[T(\mathbf{X})] \geq i_X(\theta)^{-1}[\partial g(\theta)/\partial\theta]^2, \quad (8.13)$$

where $g(\theta) := \mathbb{E}_\theta[T(\mathbf{X})]$ and

$$i_X(\theta) = \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial\theta} \log f(\mathbf{X}, \theta) \right)^2 \right] \quad (8.14)$$

is Fisher's information of $f(\mathbf{x}, \theta)$. In particular, if $\mathbb{E}_\theta[T(\mathbf{X})] = \theta$, i.e. $T(\mathbf{X})$ is an unbiased estimator of θ , then

$$\text{Var}_\theta[T(\mathbf{X})] \geq i_X(\theta)^{-1}.$$

Note that, by the independence of $\mathbf{X}_1, \dots, \mathbf{X}_N$,

$$i_X(\theta) = \text{Var}_\theta \left[\frac{\partial}{\partial\theta} \log f(\mathbf{X}, \theta) \right] = \sum_{j=1}^N \text{Var}_\theta \left[\frac{\partial}{\partial\theta} \log f_j(\mathbf{X}_j, \theta) \right] = Ni(\theta),$$

where $i(\theta) = \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial\theta} \log f_j(\mathbf{X}_j, \theta) \right)^2 \right]$ is the information number of individual \mathbf{X}_j .

Proof. We have that

$$\mathbb{E}_\theta \left[\frac{\partial}{\partial\theta} \log f(\mathbf{X}, \theta) \right] = \int \frac{\frac{\partial}{\partial\theta} f(\mathbf{x}, \theta)}{f(\mathbf{x}, \theta)} f(\mathbf{x}, \theta) d\mathbf{x} = \int \frac{\partial}{\partial\theta} f(\mathbf{x}, \theta) d\mathbf{x} = \frac{\partial}{\partial\theta} \int f(\mathbf{x}, \theta) d\mathbf{x} = 0,$$

provided the derivative can be interchanged with the integral (see Remark 8.1). Then

$$\begin{aligned} \text{Cov}_\theta \left(T(\mathbf{X}), \frac{\partial}{\partial\theta} \log f(\mathbf{X}, \theta) \right) &= \mathbb{E}_\theta \left[T(\mathbf{X}) \frac{\partial}{\partial\theta} \log f(\mathbf{X}, \theta) \right] = \mathbb{E}_\theta \left[T(\mathbf{X}) \frac{\partial}{\partial\theta} f(\mathbf{X}, \theta) / f(\mathbf{X}, \theta) \right] \\ &= \int T(\mathbf{x}) \partial f(\mathbf{x}, \theta) / \partial\theta d\mathbf{x} = \frac{\partial}{\partial\theta} \int T(\mathbf{x}) f(\mathbf{x}, \theta) d\mathbf{x}. \end{aligned}$$

That is,

$$\text{Cov}_\theta (T(\mathbf{X}), \frac{\partial}{\partial\theta} \log f(\mathbf{X}, \theta)) = \frac{\partial}{\partial\theta} \mathbb{E}_\theta[T(\mathbf{X})] = \partial g(\theta) / \partial\theta.$$

Now by Cauchy inequality we have

$$\left[\text{Cov}_\theta (T(\mathbf{X}), \frac{\partial}{\partial\theta} \log f(\mathbf{X}, \theta)) \right]^2 \leq \text{Var}_\theta[T(\mathbf{X})] \text{Var}_\theta \left[\frac{\partial}{\partial\theta} \log f(\mathbf{X}, \theta) \right].$$

Moreover

$$\text{Var}_\theta \left[\frac{\partial}{\partial\theta} \log f(\mathbf{X}, \theta) \right] = \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial\theta} \log f(\mathbf{X}, \theta) \right)^2 \right] = i(\theta),$$

and hence the inequality (8.13) follows. \square

This bound can be extended to a multivariate setting.

Theorem 8.2 (multivariate Cramér - Rao lower bound) Let $\mathbf{X}_1, \dots, \mathbf{X}_N$ be an iid sample from $f(\mathbf{x}, \boldsymbol{\theta})$, $\boldsymbol{\theta} \in \mathbb{R}^k$, and $T = T(\mathbf{X})$ be an unbiased estimator of $\boldsymbol{\theta}$, i.e., $\mathbb{E}_\theta[T(\mathbf{X})] = \boldsymbol{\theta}$. Suppose that the information matrix $\mathbf{I}(\boldsymbol{\theta})$ is nonsingular and the interchangeability property holds. Then

$$\text{Var}_\theta\left(\sum_{i=1}^k a_i T_i\right) \geq \mathbf{a}' \mathbf{I}(\boldsymbol{\theta})^{-1} \mathbf{a} \quad (8.15)$$

for any $\mathbf{a} \in \mathbb{R}^k$.

Proof. For $\mathbf{a}, \mathbf{b} \in \mathbb{R}^k$ we have (by using the interchangeability property)

$$\begin{aligned} \text{Cov}_\theta\left(\sum_{i=1}^k a_i T_i, \sum_{j=1}^k b_j \frac{\partial}{\partial \theta_j} \log f(\mathbf{X}, \boldsymbol{\theta})\right) &= \text{Cov}_\theta\left(\sum_{i=1}^k a_i T_i, \sum_{j=1}^k b_j \frac{\partial f(\mathbf{X}, \boldsymbol{\theta}) / \partial \theta_j}{f(\mathbf{X}, \boldsymbol{\theta})}\right) = \\ &= \int \left(\sum_{i=1}^k a_i T_i(\mathbf{x})\right) \left(\sum_{j=1}^k b_j \partial f(\mathbf{x}, \boldsymbol{\theta}) / \partial \theta_j\right) d\mathbf{x} = \sum_{j=1}^k b_j \frac{\partial}{\partial \theta_j} \int \left(\sum_{i=1}^k a_i T_i(\mathbf{x})\right) f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} = \\ &= \sum_{j=1}^k b_j \frac{\partial}{\partial \theta_j} \mathbb{E}_\theta \left[\sum_{i=1}^k a_i T_i\right]. \end{aligned}$$

It follows by Cauchy inequality

$$\left(\sum_{j=1}^k b_j \frac{\partial}{\partial \theta_j} \mathbb{E}_\theta \left[\sum_{i=1}^k a_i T_i\right]\right)^2 \leq \text{Var}_\theta\left(\sum_{i=1}^k a_i T_i\right) \mathbb{E}_\theta \left[\left(\sum_{j=1}^k b_j \frac{\partial}{\partial \theta_j} \log f(\mathbf{X}, \boldsymbol{\theta})\right)^2\right].$$

Since T is unbiased we have that $\mathbb{E}_\theta \left[\sum_{i=1}^k a_i T_i\right] = \sum_{i=1}^k a_i \theta_i$, and hence $\frac{\partial}{\partial \theta_j} \mathbb{E}_\theta \left[\sum_{i=1}^k a_i T_i\right] = a_j$. Also $\text{Var}_\theta\left(\sum_{i=1}^k a_i T_i\right) = \mathbf{a}' \boldsymbol{\Sigma} \mathbf{a}$, where $\boldsymbol{\Sigma}$ is the covariance matrix of T , and

$$\mathbb{E}_\theta \left[\left(\sum_{j=1}^k b_j \frac{\partial}{\partial \theta_j} \log f(\mathbf{X}, \boldsymbol{\theta})\right)^2\right] = \mathbf{b}' \mathbf{I}(\boldsymbol{\theta}) \mathbf{b}.$$

We obtain that

$$(\mathbf{a}' \mathbf{b})^2 \leq (\mathbf{a}' \boldsymbol{\Sigma} \mathbf{a})(\mathbf{b}' \mathbf{I}(\boldsymbol{\theta}) \mathbf{b}).$$

It follows that

$$\mathbf{a}' \boldsymbol{\Sigma} \mathbf{a} \geq \max_{\mathbf{b} \neq \mathbf{0}} \frac{\mathbf{b}' (\mathbf{a} \mathbf{a}') \mathbf{b}}{\mathbf{b}' \mathbf{I}(\boldsymbol{\theta}) \mathbf{b}}.$$

The maximum in the right hand side of the above inequality is attained for $\mathbf{b} = \mathbf{I}(\boldsymbol{\theta})^{-1} \mathbf{a}$ (see section 13.1.1), and hence this maximum is equal to $\mathbf{a}' \mathbf{I}(\boldsymbol{\theta})^{-1} \mathbf{a}$. Therefore we obtain that

$$\mathbf{a}' \boldsymbol{\Sigma} \mathbf{a} \geq \mathbf{a}' \mathbf{I}(\boldsymbol{\theta})^{-1} \mathbf{a} \quad (8.16)$$

for any $\mathbf{a} \in \mathbb{R}^k$. □

Definition 8.1 It is said that a sequence of estimators W_N is asymptotically efficient for $\boldsymbol{\theta}$ if $\sqrt{N}(W_N - \boldsymbol{\theta})$ converges in distribution to normal $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ with covariance matrix $\boldsymbol{\Sigma} = \mathbf{I}(\boldsymbol{\theta})^{-1}$.

The basic result (8.12), of asymptotic normality of the ML estimator $\hat{\boldsymbol{\theta}}_N$, shows that the MLE is asymptotically efficient. That is, in the standard case, under the corresponding regularity conditions, the MLE attains asymptotically the smallest possible variance. It could be noted that the bound (8.15) is not asymptotic. On the other hand, it assumes that the estimator T is unbiased, while the ML estimators often are biased. There are some other concepts of the “best possible” estimators. In the next section we briefly discuss some basic concepts.

8.3 Best unbiased estimators

Let $\mathbf{X}_1, \dots, \mathbf{X}_N$ be an iid random sample and $f(\mathbf{x}, \boldsymbol{\theta})$ be pdf of $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_N)$. By writing P_θ it is emphasized that the probability distribution of \mathbf{X} depends on the parameter vector $\boldsymbol{\theta}$.

Definition 8.2 A statistic $T(\mathbf{X})$ is a sufficient statistic for $\boldsymbol{\theta}$ if the conditional distribution of sample \mathbf{X} given $T(\mathbf{X})$ does not depend on $\boldsymbol{\theta}$. That is, $\text{Prob}(\mathbf{X} \in A | T = t)$ is independent of $\boldsymbol{\theta}$ for all (measurable) sets A and t in the range of T .

Note that a sufficient statistic always exists, take for example $T(\mathbf{X}) = \mathbf{X}$.

Theorem 8.3 (Fisher - Neyman factorization criterion) Suppose that \mathbf{X} has pdf $f(\mathbf{x}, \boldsymbol{\theta})$, $\boldsymbol{\theta} \in \Theta$. Then $T = T(\mathbf{X})$ is sufficient for $\boldsymbol{\theta}$ iff $f(\mathbf{x}, \boldsymbol{\theta}) = g(T(\mathbf{x}), \boldsymbol{\theta})h(\mathbf{x})$.

Proof. (Sketch for discrete distribution)

Suppose that T is sufficient. Then Since $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_N)$ has discrete distribution, the weight function $f(\mathbf{x}, \boldsymbol{\theta}) = P_\theta(\mathbf{X} = \mathbf{x})$, where \mathbf{x} is in the range of \mathbf{X} . Moreover

$$P_\theta(\mathbf{X} = \mathbf{x}) = \sum_t P_\theta(\mathbf{X} = \mathbf{x}, T = t),$$

where the summation over possible values of T . Also since $T = T(\mathbf{X})$,

$$\sum_t P_\theta(\mathbf{X} = \mathbf{x}, T = t) = P_\theta(\mathbf{X} = \mathbf{x}, T = T(\mathbf{x})).$$

Then

$$P_\theta(\mathbf{X} = \mathbf{x}, T = T(\mathbf{x})) = P_\theta(T = T(\mathbf{x}))P_\theta(\mathbf{X} = \mathbf{x} | T = T(\mathbf{x})),$$

where is used formula $\text{Prob}(A \cap B) = \text{Prob}(B)\text{Prob}(A|B)$ for events $B := \{T = T(\mathbf{x})\}$ and $A := \{\mathbf{X} = \mathbf{x}\}$.

By sufficiency of T we have that the conditional probability $h(\mathbf{x}) := P_\theta(\mathbf{X} = \mathbf{x} | T = T(\mathbf{x}))$ does not depend on $\boldsymbol{\theta}$. Define $g(T(\mathbf{x}), \boldsymbol{\theta}) := P_\theta(T = T(\mathbf{x}))$. It follows that $f(\mathbf{x}, \boldsymbol{\theta}) = g(T(\mathbf{x}), \boldsymbol{\theta})h(\mathbf{x})$. This shows that sufficiency implies factorization.

Now suppose that $f(\mathbf{x}, \boldsymbol{\theta}) = g(T(\mathbf{x}), \boldsymbol{\theta})h(\mathbf{x})$. Then when $T(\mathbf{x}) = t$ we have

$$\begin{aligned} P_\theta(\mathbf{X} = \mathbf{x} | T = t) &= \frac{P_\theta(\mathbf{X} = \mathbf{x}, T = t)}{P_\theta(T = t)} = \frac{g(T(\mathbf{x}), \boldsymbol{\theta})h(\mathbf{x})}{\sum_{T(\mathbf{y})=t} g(T(\mathbf{y}), \boldsymbol{\theta})h(\mathbf{y})} \\ &= \frac{g(T(\mathbf{x}), \boldsymbol{\theta})h(\mathbf{x})}{\sum_{T(\mathbf{y})=t} f(\mathbf{y}, \boldsymbol{\theta})} = \frac{g(t, \boldsymbol{\theta})h(\mathbf{x})}{\sum_{T(\mathbf{y})=t} g(t, \boldsymbol{\theta})h(\mathbf{y})} = \frac{h(\mathbf{x})}{\sum_{T(\mathbf{y})=t} h(\mathbf{y})}, \end{aligned}$$

which does not depend on $\boldsymbol{\theta}$. If $T(\mathbf{x}) \neq t$, then $P_\theta(\mathbf{X} = \mathbf{x} | T = t) = 0$. It follows that $T(\mathbf{X})$ is sufficient. \square

Let $T(\mathbf{X})$ be a sufficient statistic for $\boldsymbol{\theta}$. Then by the Factorization Theorem, the likelihood function

$$L_N(\boldsymbol{\theta}) = f(\mathbf{x}, \boldsymbol{\theta}) = g(T(\mathbf{x}), \boldsymbol{\theta})h(\mathbf{x}).$$

It follows that the MLE $\hat{\boldsymbol{\theta}}$ is a function of T , i.e.,

$$\hat{\boldsymbol{\theta}} \in \arg \max_{\boldsymbol{\theta}} g(T(\mathbf{x}), \boldsymbol{\theta}).$$

Definition 8.3 A sufficient statistic $T = T(\mathbf{X})$ is said to be minimal sufficient if for any other sufficient statistic $S = S(\mathbf{X})$, there exists a function $g(\cdot)$ such that $T = g(S)$.

Theorem 8.4 (Lehmann - Scheffe) Suppose that there exists $T(\mathbf{X})$ such that for any \mathbf{x} and \mathbf{y} , the ratio $\frac{f(\mathbf{x}, \boldsymbol{\theta})}{f(\mathbf{y}, \boldsymbol{\theta})}$ is independent of $\boldsymbol{\theta}$ if and only if $T(\mathbf{x}) = T(\mathbf{y})$. Then $T(\mathbf{X})$ is a minimal sufficient statistic for $\boldsymbol{\theta}$.

Proof. Let us show that $T(\mathbf{X})$ is sufficient. For t in the image of $T(\mathbf{x})$ consider sets $A_t = \{\mathbf{x} : T(\mathbf{x}) = t\}$. For t in the image of $T(\mathbf{x})$, consider a point $\mathbf{x}_t \in A_t$. We have that $\mathbf{x}_{T(\mathbf{x})}$ and \mathbf{x} are in the same set A_t , i.e., $T(\mathbf{x}) = T(\mathbf{x}_{T(\mathbf{x})})$, and hence by the assumption of the theorem, the ratio $f(\mathbf{x}, \boldsymbol{\theta})/f(\mathbf{x}_{T(\mathbf{x})}, \boldsymbol{\theta})$ does not depend on $\boldsymbol{\theta}$. Define $h(\mathbf{x}) = f(\mathbf{x}, \boldsymbol{\theta})/f(\mathbf{x}_{T(\mathbf{x})}, \boldsymbol{\theta})$ and $g(t, \boldsymbol{\theta}) = f(\mathbf{x}_t, \boldsymbol{\theta})$. Then

$$f(\mathbf{x}, \boldsymbol{\theta}) = \frac{f(\mathbf{x}_{T(\mathbf{x})}, \boldsymbol{\theta})f(\mathbf{x}, \boldsymbol{\theta})}{f(\mathbf{x}_{T(\mathbf{x})}, \boldsymbol{\theta})} = g(T(\mathbf{x}), \boldsymbol{\theta})h(\mathbf{x}).$$

It follows by the Factorization Theorem that $T(\mathbf{X})$ is a sufficient statistic.

Let us show that $T(\mathbf{X})$ is minimal sufficient. Let $T'(\mathbf{X})$ be a sufficient statistic. By the Factorization Theorem, $f(\mathbf{x}, \boldsymbol{\theta}) = g(T'(\mathbf{x}), \boldsymbol{\theta})h(\mathbf{x})$. Suppose that $T'(\mathbf{x}) = T'(\mathbf{y})$. Then

$$\frac{f(\mathbf{x}, \boldsymbol{\theta})}{f(\mathbf{y}, \boldsymbol{\theta})} = \frac{g(T'(\mathbf{x}), \boldsymbol{\theta})h(\mathbf{x})}{g(T'(\mathbf{y}), \boldsymbol{\theta})h(\mathbf{y})} = \frac{h(\mathbf{x})}{h(\mathbf{y})}.$$

Since this ratio does not depend on $\boldsymbol{\theta}$, it follows by the assumption of theorem that $T(\mathbf{x}) = T(\mathbf{y})$. That is, $T'(\mathbf{x}) = T'(\mathbf{y})$ implies that $T(\mathbf{x}) = T(\mathbf{y})$. It follows that $T(\mathbf{X})$ is a function of $T'(\mathbf{X})$. \square

Note that the second part of the above proof shows that a sufficient statistic $T(\mathbf{X})$ is minimal sufficient if the following implication holds: if the ratio $f(\mathbf{x}, \boldsymbol{\theta})/f(\mathbf{y}, \boldsymbol{\theta})$ does not depend on $\boldsymbol{\theta}$, then $T(\mathbf{x}) = T(\mathbf{y})$.

Example 8.2 Consider exponential family of distributions in the canonical form (see eq. (7.1)),

$$f(\mathbf{x}, \boldsymbol{\theta}) = \exp \left\{ \sum_{i=1}^k \theta_i T_i(\mathbf{x}) - A(\boldsymbol{\theta}) \right\} h(\mathbf{x}), \quad (8.17)$$

with parameter space

$$\Theta = \left\{ \boldsymbol{\theta} : \int \exp \left\{ \sum_{i=1}^k \theta_i T_i(\mathbf{x}) \right\} h(\mathbf{x}) d\mathbf{x} < \infty \right\}.$$

It follows by the Factorization Theorem that $(T_1(\mathbf{X}), \dots, T_k(\mathbf{X}))$ is a sufficient statistic. Note that the set Θ is convex. Also

$$\frac{f(\mathbf{x}, \boldsymbol{\theta})}{f(\mathbf{y}, \boldsymbol{\theta})} = \exp \left\{ \sum_{i=1}^k \theta_i (T_i(\mathbf{x}) - T_i(\mathbf{y})) \right\} \frac{h(\mathbf{x})}{h(\mathbf{y})}.$$

Suppose that the set Θ has a nonempty interior. Then if the ratio $f(\mathbf{x}, \boldsymbol{\theta})/f(\mathbf{y}, \boldsymbol{\theta})$ does not depend on $\boldsymbol{\theta}$, then $T_i(\mathbf{x}) = T_i(\mathbf{y})$, $i = 1, \dots, k$. Indeed if this ratio does not depend on $\boldsymbol{\theta}$, then

$$\frac{\partial}{\partial \theta_i} \exp \left\{ \sum_{i=1}^k \theta_i (T_i(\mathbf{x}) - T_i(\mathbf{y})) \right\} \frac{h(\mathbf{x})}{h(\mathbf{y})} = (T_i(\mathbf{x}) - T_i(\mathbf{y})) \exp \left\{ \sum_{i=1}^k \theta_i (T_i(\mathbf{x}) - T_i(\mathbf{y})) \right\} \frac{h(\mathbf{x})}{h(\mathbf{y})}$$

is zero at every interior point of the set Θ . It follows that $T_i(\mathbf{x}) = T_i(\mathbf{y})$. This implies that $(T_1(\mathbf{X}), \dots, T_k(\mathbf{X}))$ is minimal sufficient. \square

We assume in the remainder of this section that $g : \Theta \rightarrow \mathbb{R}$ is a real valued (measurable) function.

Definition 8.4 An estimator $T = T(\mathbf{X})$ of $g(\boldsymbol{\theta})$ is a best unbiased estimator if $\mathbb{E}_\theta[T] = g(\boldsymbol{\theta})$ for all $\boldsymbol{\theta} \in \Theta$, and for any unbiased estimator $S = S(\mathbf{X})$ of $g(\boldsymbol{\theta})$ it follows that

$$\text{Var}_\theta[T] \leq \text{Var}_\theta[S], \quad \forall \boldsymbol{\theta} \in \Theta.$$

Best unbiased estimator is called Uniform Minimum Variance Unbiased (UMVU) estimator.

Finding an UMVU estimator could be not easy. The following result shows that conditioning of any unbiased estimator on a sufficient statistic will result in uniform reduction of the variance. Therefore if an UMVU estimator exists, then it is a function of (minimal) sufficient statistic.

For random variables X and Y , we use below property $\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|Y]]$ of conditional expectation, and the following formula

$$\text{Var}(X) = \mathbb{E}[\text{Var}(X|Y)] + \text{Var}[\mathbb{E}(X|Y)], \quad (8.18)$$

for conditional variance

$$\text{Var}(X|Y) = \mathbb{E}[(X - \mathbb{E}(X|Y))^2|Y].$$

Indeed

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}[(X - \mathbb{E}(X))^2] = \mathbb{E}[\mathbb{E}[(X - \mathbb{E}(X))^2|Y]] \\ &= \mathbb{E} \left[\mathbb{E}[(X - \mathbb{E}(X|Y) + \mathbb{E}(X|Y) - \mathbb{E}(X))^2|Y] \right] \\ &= \mathbb{E} \left[\underbrace{\mathbb{E}[(X - \mathbb{E}(X|Y))^2|Y]}_{\text{Var}(X|Y)} + \underbrace{\mathbb{E}[(\mathbb{E}(X|Y) - \mathbb{E}(X))^2]}_{\text{Var}[\mathbb{E}(X|Y)]} \right]. \end{aligned}$$

In the above derivation we used that

$$\mathbb{E}[(\mathbb{E}(X|Y) - \mathbb{E}(X))^2|Y] = \mathbb{E}[(\mathbb{E}(X|Y) - \mathbb{E}(X))^2],$$

since $(\mathbb{E}(X|Y) - \mathbb{E}(X))^2$ is a function of Y , and that

$$\mathbb{E}[(\mathbb{E}[(X - \mathbb{E}(X|Y))](\mathbb{E}(X|Y) - \mathbb{E}(X))|Y] = \mathbb{E}[(\mathbb{E}(X - \mathbb{E}(X|Y)|Y)(\mathbb{E}(X|Y) - \mathbb{E}(X)))] = 0.$$

Theorem 8.5 (Rao - Blackwell) Let W be an unbiased estimator of $g(\boldsymbol{\theta})$, and T be a sufficient statistic for $\boldsymbol{\theta}$. Define $h(t) := \mathbb{E}[W|T = t]$. Then $\mathbb{E}_\theta[h(T)] = g(\boldsymbol{\theta})$ and

$$\text{Var}_\theta[h(T)] \leq \text{Var}_\theta[W], \quad \forall \boldsymbol{\theta} \in \Theta. \quad (8.19)$$

Moreover, unless $P_\theta\{W = h(T)\} = 1$, the inequality (8.19) is strict.

Proof. Note that by sufficiency of T , $h(T)$ does not depend on $\boldsymbol{\theta}$ and hence is a statistic. We have that

$$\mathbb{E}_\theta[h(T)] = \mathbb{E}_\theta[\mathbb{E}[W|T = t]] = \mathbb{E}_\theta[W] = g(\boldsymbol{\theta}).$$

That is $h(T)$ is an unbiased estimator of $g(\boldsymbol{\theta})$. Now by using formula (8.18) for conditional variance

$$\begin{aligned} \text{Var}_\theta[W] &= \text{Var}_\theta[\mathbb{E}(W|T)] + \mathbb{E}_\theta[\text{Var}(W|T)] \\ &= \text{Var}_\theta[h(T)] + \mathbb{E}_\theta[\text{Var}(W|T)] \geq \text{Var}_\theta[h(T)], \end{aligned}$$

and hence (8.19) follows. Moreover, $\mathbb{E}_\theta[\text{Var}(W|T)] > 0$ and hence the inequality (8.19) is strict unless $P_\theta\{W = h(T)\} = 1$. \square

Theorem 8.6 An UMVU estimator W (if it exists) of $g(\boldsymbol{\theta})$ is unique.

Proof. Let W' be another UMVU estimator of $g(\boldsymbol{\theta})$. Then $W^* = (W + W')/2$ is unbiased and

$$\text{Var}_\theta(W^*) = \frac{1}{4}\text{Var}_\theta(W) + \frac{1}{4}\text{Var}_\theta(W') + \frac{1}{2}\text{Cov}_\theta(W, W').$$

Now by Cauchy inequality

$$\text{Cov}_\theta(W, W') \leq [\text{Var}_\theta(W) \cdot \text{Var}_\theta(W')]^{1/2},$$

and $\text{Var}_\theta(W) = \text{Var}_\theta(W')$ by the minimum variance assumption. Hence $\text{Var}_\theta(W^*) \leq \text{Var}_\theta(W)$. Since W is UMVU it follows that $\text{Var}_\theta(W^*) = \text{Var}_\theta(W)$ for all $\boldsymbol{\theta} \in \Theta$. The equality in Cauchy inequality holds only if $W' = a(\boldsymbol{\theta})W + b(\boldsymbol{\theta})$. Then

$$\text{Cov}_\theta(W, W') = \text{Cov}_\theta(W, a(\boldsymbol{\theta})W + b(\boldsymbol{\theta})) = a(\boldsymbol{\theta})\text{Var}_\theta(W).$$

Also by the above we have that $\text{Cov}_\theta(W, W') = \text{Var}_\theta(W)$ and hence $a(\boldsymbol{\theta}) \equiv 1$. Moreover $\mathbb{E}_\theta[W'] = g(\boldsymbol{\theta}) = \mathbb{E}_\theta[W]$ and hence $b(\boldsymbol{\theta}) \equiv 0$. It follows that $W = W'$. \square

Definition 8.5 Loss function (for estimating $g(\boldsymbol{\theta})$) is a nonnegative valued function $L(\boldsymbol{\theta}, a)$, $\boldsymbol{\theta} \in \Theta$, $a \in \mathbb{R}$, such that $L(\boldsymbol{\theta}, g(\boldsymbol{\theta})) = 0$ for all $\boldsymbol{\theta} \in \Theta$. Risk function $R(\boldsymbol{\theta}, T) := \mathbb{E}_\theta[L(\boldsymbol{\theta}, T(\mathbf{X}))]$, where $T(\mathbf{X})$ is an estimator of $g(\boldsymbol{\theta})$.

For example $L(\boldsymbol{\theta}, a) := |a - g(\boldsymbol{\theta})|^p$, $p > 0$, is a loss function. If $L(\boldsymbol{\theta}, a) = (g(\boldsymbol{\theta}) - a)^2$, then

$$R(\boldsymbol{\theta}, T) = \mathbb{E}_\theta[(g(\boldsymbol{\theta}) - T(\mathbf{X}))^2]$$

is the Mean Square Error of estimator T of $g(\boldsymbol{\theta})$.

Theorem 8.7 (Another version of Rao - Blackwell theorem) Let $L(\boldsymbol{\theta}, a)$ be a loss function, W be a sufficient statistic and $h(t) = \mathbb{E}[W|T = t]$. Suppose that $L(\boldsymbol{\theta}, \cdot)$ is strictly convex. Then

$$R(\boldsymbol{\theta}, h(T)) \leq R(\boldsymbol{\theta}, W), \tag{8.20}$$

and the above inequality is strict unless $P_\theta\{W = h(T)\} = 1$.

Proof. By using Jensen's inequality

$$R(\boldsymbol{\theta}, h(T)) = \mathbb{E}_\theta[L(\boldsymbol{\theta}, \mathbb{E}[W|T])] \leq \mathbb{E}_\theta[\mathbb{E}[L(\boldsymbol{\theta}, W)|T]] = \mathbb{E}_\theta[L(\boldsymbol{\theta}, W)] = R(\boldsymbol{\theta}, W).$$

The inequality (8.20) follows and this inequality is strict unless $P_\theta\{W = h(T)\} = 1$. \square

9 Hypotheses testing

Let $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_N)$ be a random sample (data). Consider testing $H_0 : \boldsymbol{\theta} \in \Theta_0$ versus $H_1 : \boldsymbol{\theta} \in \Theta_1$, where $\Theta_0, \Theta_1 \subset \mathbb{R}^k$. A procedure for such testing consists of choosing a set $R \subset \mathbb{R}^d$, referred to as the rejection region, and hence defining its complement $R^c = \mathbb{R}^d \setminus R$ referred to as the acceptance region, where d is the dimension of \mathbf{X} . That is, reject H_0 if $\mathbf{X} \in R$. Alternatively this can be formulated as accept H_0 if $\mathbf{X} \in R^c$. Rejecting H_0 automatically means acceptance H_1 , and acceptance H_0 means rejection of H_1 .

There are two types of errors, type I error - reject H_0 when H_0 is true, type II error - accept H_0 when H_0 is false. The corresponding probabilities $\alpha = P_\theta(\text{type I error})$ and $\beta = P_\theta(\text{type II error})$. That is

$$\begin{aligned}\alpha &= P_\theta(\mathbf{X} \in R), \quad \boldsymbol{\theta} \in \Theta_0, \\ \beta &= P_\theta(\mathbf{X} \in R^c), \quad \boldsymbol{\theta} \in \Theta_1.\end{aligned}$$

Power of the test is $1 - \beta = P_\theta(\mathbf{X} \in R)$, $\boldsymbol{\theta} \in \Theta_1$. Note that $\alpha = \alpha(\boldsymbol{\theta})$ and $\beta = \beta(\boldsymbol{\theta})$ are functions of $\boldsymbol{\theta}$.

Theorem 9.1 (Neyman - Pearson Lemma) *Consider simple alternatives $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ versus $H_1 : \boldsymbol{\theta} = \boldsymbol{\theta}_1$ with respective pdfs $f(\mathbf{x}, \boldsymbol{\theta}_0)$ and $f(\mathbf{x}, \boldsymbol{\theta}_1)$. Then the minimal error rejection region is*

$$R = \{\mathbf{x} \in \mathbb{R}^d : f(\mathbf{x}, \boldsymbol{\theta}_1) \geq \kappa f(\mathbf{x}, \boldsymbol{\theta}_0)\}, \quad (9.1)$$

where $\kappa > 0$ is such that $\int_R f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x} = \alpha$.

Proof. Note that $\int_R f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x} = \alpha$ and $\int_{R^c} f(\mathbf{x}, \boldsymbol{\theta}_1) d\mathbf{x} = \beta$. We want to choose region R , or equivalently R^c , such that the probability of type I error equals the significance level α , and the probability of type II error is the smallest possible. For a constant $\kappa > 0$ this can be formulated as minimization of

$$\int_{R^c} f(\mathbf{x}, \boldsymbol{\theta}_1) d\mathbf{x} + \kappa \int_R f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x},$$

with respect to R (or equivalently with respect to R^c), subject to $\int_R f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x} = \alpha$. Since $R = \mathbb{R}^d \setminus R^c$ we have that

$$\int_R f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x} = \int_{\mathbb{R}^d} f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x} - \int_{R^c} f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x}.$$

Moreover $\int_{\mathbb{R}^d} f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x} = 1$, and hence

$$\int_{R^c} f(\mathbf{x}, \boldsymbol{\theta}_1) d\mathbf{x} + \kappa \int_R f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x} = \kappa + \int_{R^c} [f(\mathbf{x}, \boldsymbol{\theta}_1) - \kappa f(\mathbf{x}, \boldsymbol{\theta}_0)] d\mathbf{x}.$$

It follows that the minimum is attained for

$$R^c = \{\mathbf{x} : f(\mathbf{x}, \boldsymbol{\theta}_1) - \kappa f(\mathbf{x}, \boldsymbol{\theta}_0) < 0\},$$

or equivalently for

$$R = \{\mathbf{x} : f(\mathbf{x}, \boldsymbol{\theta}_1) - \kappa f(\mathbf{x}, \boldsymbol{\theta}_0) \geq 0\}.$$

Note that for $\kappa = 0$ the rejection region $R = \mathbb{R}^d$ and hence $\alpha = 1$. By increasing κ the rejection region shrinks and $\int_R f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x}$ continuously decreases and tends to zero. Therefore we can choose κ such that $\int_R f(\mathbf{x}, \boldsymbol{\theta}_0) d\mathbf{x} = \alpha$. \square

Suppose that $T(\mathbf{X})$ is a sufficient statistic for $\boldsymbol{\theta}$. By Factorization Theorem (Theorem 8.3), $f(\mathbf{x}, \boldsymbol{\theta}_0) = g(T(\mathbf{x}), \boldsymbol{\theta}_0)h(\mathbf{x})$ and $f(\mathbf{x}, \boldsymbol{\theta}_1) = g(T(\mathbf{x}), \boldsymbol{\theta}_1)h(\mathbf{x})$. Therefore the rejection region (9.1) can be written as

$$R = \{\mathbf{x} : g(T(\mathbf{x}), \boldsymbol{\theta}_1) \geq \kappa g(T(\mathbf{x}), \boldsymbol{\theta}_0)\}.$$

9.1 Likelihood Ratio Test

Consider

$$\lambda(\mathbf{x}) := \frac{\sup_{\boldsymbol{\theta} \in \Theta} L(\boldsymbol{\theta})}{\sup_{\boldsymbol{\theta} \in \Theta_0} L(\boldsymbol{\theta})},$$

where $\Theta = \Theta_0 \cup \Theta_1$ and $L(\boldsymbol{\theta}) = f(\mathbf{x}, \boldsymbol{\theta})$ is the corresponding likelihood function. Note that $\lambda(\mathbf{x}) \geq 1$ since $\Theta_0 \subset \Theta$. The rejection region of the Likelihood Ratio Test (LRT) is

$$R = \{\mathbf{x} : \lambda(\mathbf{x}) \geq c\},$$

for some $c > 1$. That is, the H_0 is rejected for large value of the LRT statistic.

If $T(\mathbf{X})$ is a sufficient statistic for $\boldsymbol{\theta}$, then by the Factorization Theorem

$$\lambda(\mathbf{x}) = \frac{\sup_{\boldsymbol{\theta} \in \Theta} g(T(\mathbf{x}, \boldsymbol{\theta}))}{\sup_{\boldsymbol{\theta} \in \Theta_0} g(T(\mathbf{x}, \boldsymbol{\theta}))}.$$

That is, the LRT can be formulated in terms of the sufficient statistic $T(\mathbf{X})$. For simple alternatives when $\Theta_0 = \{\boldsymbol{\theta}_0\}$ and $\Theta_1 = \{\boldsymbol{\theta}_1\}$ we have that

$$\lambda(\mathbf{x}) = \frac{\max\{L(\boldsymbol{\theta}_0), L(\boldsymbol{\theta}_1)\}}{L(\boldsymbol{\theta}_0)} = \max\{1, f(\mathbf{x}, \boldsymbol{\theta}_1)/f(\mathbf{x}, \boldsymbol{\theta}_0)\},$$

and hence this is equivalent to the rejection region of the Neyman - Pearson Lemma.

Let us discuss asymptotics of the LRT. We will discuss this for the simple hypothesis $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ against the unrestricted alternative $H_1 : \boldsymbol{\theta} \in \mathbb{R}^k$. We have that

$$2 \log \lambda(\mathbf{X}) = -2 \log L(\boldsymbol{\theta}_0) + 2 \sup_{\boldsymbol{\theta} \in \mathbb{R}^k} \log L(\boldsymbol{\theta}).$$

Note that

$$\sup_{\boldsymbol{\theta} \in \mathbb{R}^k} \log L(\boldsymbol{\theta}) = \log L(\hat{\boldsymbol{\theta}}),$$

where $\hat{\boldsymbol{\theta}}$ is the ML estimator under the unrestricted alternative H_1 . Consider

$$S(\boldsymbol{\theta}) := \frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}) = \sum_{i=1}^N \frac{\partial}{\partial \boldsymbol{\theta}} \log f(X_i, \boldsymbol{\theta}), \quad (9.2)$$

called the *score function*. Note that $S(\hat{\boldsymbol{\theta}}) = \mathbf{0}$ (necessary optimality condition), and $\mathbb{E}_{\boldsymbol{\theta}}[S(\boldsymbol{\theta})] = \mathbf{0}$ (see equation (8.7)). Now using second order Taylor approximation,

$$\log L(\hat{\boldsymbol{\theta}}) \approx \log L(\boldsymbol{\theta}_0) + \left[\frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}_0) \right]' (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + \frac{1}{2} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \left[\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log L(\boldsymbol{\theta}_0) \right] (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0).$$

Note that $\frac{\partial}{\partial \boldsymbol{\theta}} \log L(\boldsymbol{\theta}_0) = S(\boldsymbol{\theta}_0)$ and $\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log L(\boldsymbol{\theta}_0) = \frac{\partial}{\partial \boldsymbol{\theta}'} S(\boldsymbol{\theta}_0)$. Hence and since $S(\hat{\boldsymbol{\theta}}) = \mathbf{0}$,

$$\log L(\hat{\boldsymbol{\theta}}) \approx \log L(\boldsymbol{\theta}_0) - [S(\hat{\boldsymbol{\theta}}) - S(\boldsymbol{\theta}_0)]' (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) + \frac{1}{2} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \left[\frac{\partial}{\partial \boldsymbol{\theta}'} S(\boldsymbol{\theta}_0) \right] (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0).$$

Also first order approximation of the score function:

$$S(\hat{\boldsymbol{\theta}}) - S(\boldsymbol{\theta}_0) \approx \left[\frac{\partial}{\partial \boldsymbol{\theta}'} S(\boldsymbol{\theta}_0) \right] (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0).$$

Therefore

$$\begin{aligned}
2 \log \lambda(\mathbf{X}) &= -2 \log L(\boldsymbol{\theta}_0) + 2 \log L(\hat{\boldsymbol{\theta}}) \\
&\approx (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \left[-\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log L(\boldsymbol{\theta}_0) \right] (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \\
&= [\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)]' \left[-\frac{1}{N} \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log L(\boldsymbol{\theta}_0) \right] [\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)].
\end{aligned}$$

Assuming H_0 , we have that $\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)$ converges in distribution to $\mathcal{N}(\mathbf{0}, I(\boldsymbol{\theta}_0)^{-1})$, and by (8.4) and the LLN,

$$-\frac{1}{N} \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log L(\boldsymbol{\theta}_0) = -\frac{1}{N} \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \sum_{i=1}^N \log f(\mathbf{X}_i, \boldsymbol{\theta}_0)$$

converges in probability to $I(\boldsymbol{\theta}_0)$. It follows that under H_0 , the statistic $2 \log \lambda(\mathbf{X})$ converges in distribution to the quadratic form $Z' [I(\boldsymbol{\theta}_0)] Z$, where $Z \sim \mathcal{N}(\mathbf{0}, I(\boldsymbol{\theta}_0)^{-1})$. By theorem 3.1 this implies that $2 \log \lambda(\mathbf{X})$ converges in distribution to χ_k^2 . \square

In general $2 \log \lambda(\mathbf{X})$ converges in distribution to χ_{k-q}^2 under H_0 , where $\Theta_0 \subset \mathbb{R}^k$ is a smooth manifold of dimension $q = \dim \Theta_0$.

Power of the LRT under local alternatives

Suppose the following so-called parameter drift (local alternatives) for testing $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ against $H_1 : \boldsymbol{\theta}_{0,N} = \boldsymbol{\theta}_0 + N^{-1/2} \mathbf{b}$, where $\mathbf{b} \in \mathbb{R}^k$ is a fixed vector. Then

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = \sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_{0,N}) + \sqrt{N}(\boldsymbol{\theta}_{0,N} - \boldsymbol{\theta}_0) = \sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_{0,N}) + \mathbf{b} \rightsquigarrow N(\mathbf{b}, I(\boldsymbol{\theta}_0)^{-1}).$$

Hence under local alternatives

$$2 \log \lambda \approx [\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)]' I(\boldsymbol{\theta}_0) [\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)]$$

can be approximated by the noncentral chi-square distribution $\chi_k^2(\delta)$ with the noncentrality parameter $\delta = \mathbf{b}' I(\boldsymbol{\theta}_0) \mathbf{b}$ (Theorem 3.3).

9.2 Testing equality constraints

Consider testing $H_0 : \mathbf{a}(\boldsymbol{\theta}) = (a_1(\boldsymbol{\theta}), \dots, a_q(\boldsymbol{\theta}))' = \mathbf{0}$ against $H_1 : \mathbf{a}(\boldsymbol{\theta}) \neq \mathbf{0}$. Let

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta} \in \mathbb{R}^k} L(\boldsymbol{\theta}) \text{ and } \tilde{\boldsymbol{\theta}} = \arg \max_{\mathbf{a}(\boldsymbol{\theta})=\mathbf{0}} L(\boldsymbol{\theta})$$

be the respective unrestricted and restricted ML estimators. We have here that the 2log Likelihood Ratio Test (LRT) statistic is $2[\log L(\hat{\boldsymbol{\theta}}) - \log L(\tilde{\boldsymbol{\theta}})]$. Under H_0 (and the regularity conditions) this test statistic converges in distribution to χ_q^2 .

Wald test statistic. Consider testing (linear⁴) equality constraints $H_0 : \mathbf{A}\boldsymbol{\theta} = \mathbf{c}$ against $H_1 : \mathbf{A}\boldsymbol{\theta} \neq \mathbf{c}$, where \mathbf{A} is $q \times k$ matrix of full row rank q . The Wald test statistic is

$$W := N(\mathbf{A}\hat{\boldsymbol{\theta}} - \mathbf{c})' (\mathbf{A} I(\hat{\boldsymbol{\theta}})^{-1} \mathbf{A}')^{-1} (\mathbf{A}\hat{\boldsymbol{\theta}} - \mathbf{c}).$$

⁴For nonlinear constraints we can use $\mathbf{A} = \partial \mathbf{a}(\hat{\boldsymbol{\theta}}) / \partial \boldsymbol{\theta}'$.

Suppose that the corresponding regularity conditions hold so that the (unrestricted) ML estimator $\hat{\boldsymbol{\theta}}$ is a consistent estimator of the population value $\boldsymbol{\theta}^*$ of the parameter vector, and $\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \rightsquigarrow \mathcal{N}(\mathbf{0}, I(\boldsymbol{\theta}^*)^{-1})$. Then under H_0 (i.e., $\mathbf{A}\boldsymbol{\theta}^* = \mathbf{c}$)

$$\sqrt{N}\mathbf{A}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) = \sqrt{N}(\mathbf{A}\hat{\boldsymbol{\theta}} - \mathbf{c}) \rightsquigarrow \mathcal{N}(\mathbf{0}, \mathbf{A}I(\boldsymbol{\theta}^*)^{-1}\mathbf{A}').$$

It follows that under H_0 the Wald test statistic converges in distribution to $Z'(\mathbf{A}I(\boldsymbol{\theta}^*)^{-1}\mathbf{A}')^{-1}Z$, where $Z \sim \mathcal{N}(\mathbf{0}, \mathbf{A}I(\boldsymbol{\theta}^*)^{-1}\mathbf{A}')$. Therefore by Theorem 3.1, under H_0 the Wald test statistic converges in distribution to χ_q^2 .

Note that the LRT

$$2[\log L(\hat{\boldsymbol{\theta}}) - \log L(\tilde{\boldsymbol{\theta}})] \approx \inf_{\mathbf{A}\boldsymbol{\theta}=\mathbf{c}} [\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)]' \left[-\frac{1}{N} \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log L(\boldsymbol{\theta}^*) \right] [\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)].$$

Also using formula (8.4) for the information matrix and $L(\boldsymbol{\theta}^*) = f(\mathbf{x}, \boldsymbol{\theta}^*)$, by the LLN we have that $-\frac{1}{N} \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \log L(\boldsymbol{\theta}^*) \approx I(\boldsymbol{\theta}^*)$ under H_0 . Therefore under H_0 , the LRT and Wald test statistics are asymptotically equivalent.

Score function test statistic. Consider testing $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ against $H_1 : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$. The score function test statistic is

$$N^{-1}S(\boldsymbol{\theta}_0)'I(\boldsymbol{\theta}_0)^{-1}S(\boldsymbol{\theta}_0),$$

where $S(\boldsymbol{\theta})$ is the score function (see equation (9.2)). Recall that $\mathbb{E}_{\boldsymbol{\theta}}[S(\boldsymbol{\theta})] = 0$ and $N^{-1/2}S(\boldsymbol{\theta})$ converges in distribution to $\mathcal{N}(\mathbf{0}, I(\boldsymbol{\theta}))$. It follows that under H_0 ,

$$N^{-1}S(\boldsymbol{\theta}_0)'I(\boldsymbol{\theta}_0)^{-1}S(\boldsymbol{\theta}_0) \rightsquigarrow \chi_k^2.$$

In general, under H_0 ,

$$N^{-1}S(\tilde{\boldsymbol{\theta}})'I(\tilde{\boldsymbol{\theta}})^{-1}S(\tilde{\boldsymbol{\theta}}) \rightsquigarrow \chi_q^2$$

when testing q equality constraints.

10 Multinomial distribution

Consider $\mathbf{Y} = (Y_1, \dots, Y_k)'$ with $Y_1 + \dots + Y_k = N$ and

$$\text{Prob}(\mathbf{Y} = \mathbf{y}) = \frac{N!}{y_1! \times \dots \times y_k!} \prod_{i=1}^k p_i^{y_i},$$

where $p_i > 0$, $i = 1, \dots, k$, and $p_1 + \dots + p_k = 1$. We denote this as $\mathbf{Y} \sim \text{Mult}(N, \mathbf{p})$, where $\mathbf{p} = (p_1, \dots, p_k)'$. In particular, for $k = 2$ this becomes binomial distribution $Y \sim B(N, p)$ with $\text{Prob}(Y = y) = \binom{N}{y} p^y (1-p)^{N-y}$, $y = 0, 1, \dots, N$.

The log-likelihood function, up to a constant independent of \mathbf{p} , is $L(\mathbf{p}) = \sum_{i=1}^k Y_i \log p_i$. Therefore the ML estimator of \mathbf{p} is given by the solution of the problem:

$$\max_{p_i \geq 0} \sum_{i=1}^k Y_i \log p_i \quad \text{subject to } p_1 + \dots + p_k = 1.$$

It follows that the ML estimators are $\hat{p}_i = Y_i/N$, $i = 1, \dots, k$.

If $\mathbf{Y} \sim \text{Mult}(N, \mathbf{p})$, then the covariance matrix $\text{Cov}(\mathbf{Y}) = N\mathbf{C}$, where $c_{ii} = p_i(1-p_i)$, $i = 1, \dots, k$ and $c_{ij} = -p_i p_j$, $i \neq j$.

Indeed each Y_i has binomial distribution with probability of success p_i and hence $\text{Var}(Y_i) = Np_i(1 - p_i)$. Moreover, $Y_i + Y_j$, $i \neq j$, has binomial distribution with probability of success $p_i + p_j$ and hence

$$\text{Var}(Y_i + Y_j) = N(p_i + p_j)(1 - p_i - p_j) = N(p_i - p_i^2 + p_j - p_j^2 - 2p_i p_j).$$

On the other hand

$$\text{Var}(Y_i + Y_j) = \text{Var}(Y_i) + \text{Var}(Y_j) + 2\text{Cov}(Y_i, Y_j),$$

and $\text{Var}(Y_i) = N(p_i - p_i^2)$, $\text{Var}(Y_j) = N(p_j - p_j^2)$. It follows that $\text{Cov}(Y_i, Y_j) = -Np_i p_j$.

This can be written as $\mathbf{C} = \mathbf{P} - \mathbf{p}\mathbf{p}'$, where $\mathbf{P} = \text{diag}(p_1, \dots, p_k)$ and $\mathbf{p} = (p_1, \dots, p_k)'$. Note that $\mathbf{C}\mathbf{1}_k = \mathbf{0}$ and $\text{rank}(\mathbf{C}) = k - 1$.

Consider testing $H_0 : \mathbf{p} = \mathbf{p}^*$ against $H_1 : \mathbf{p} \neq \mathbf{p}^*$. The corresponding log LRT statistic is

$$\log \lambda = \sum_{i=1}^k Y_i \log \hat{p}_i - \sum_{i=1}^k Y_i \log p_i^* = \sum_{i=1}^k Y_i \log \frac{\hat{p}_i}{p_i^*}.$$

Note that (second order Taylor approximation of $\log x$ at $x = 1$)

$$\log x = x - 1 - \frac{1}{2}(x - 1)^2 + o(x - 1)^2.$$

Under H_0 values \hat{p}_i are close to p_i^* and hence

$$\sum_{i=1}^k Y_i \log \frac{\hat{p}_i}{p_i^*} = - \sum_{i=1}^k Y_i \log \frac{p_i^*}{\hat{p}_i} \approx - \sum_{i=1}^k Y_i \left(\frac{p_i^*}{\hat{p}_i} - 1 \right) + \frac{1}{2} \sum_{i=1}^k Y_i \left(\frac{p_i^*}{\hat{p}_i} - 1 \right)^2.$$

Moreover

$$\sum_{i=1}^k Y_i \left(\frac{p_i^*}{\hat{p}_i} - 1 \right) = \sum_{i=1}^k (Np_i^* - Y_i) = 0,$$

since $\sum_{i=1}^k p_i^* = 1$ and $\sum_{i=1}^k Y_i = N$. Hence under H_0 ,

$$2 \log \lambda = 2 \sum_{i=1}^k Y_i \log \frac{\hat{p}_i}{p_i^*} \approx \sum_{i=1}^k \frac{(Y_i - Np_i^*)^2}{Y_i} \approx \sum_{i=1}^k \frac{(Y_i - Np_i^*)^2}{Np_i^*},$$

where in the last approximation we used $p_i^* \approx \hat{p}_i = Y_i/N$. Values Y_i are called observed frequencies, Np_i^* are called expected frequencies, and $\sum_{i=1}^m \frac{(Y_i - Np_i^*)^2}{Np_i^*}$ is the famous Pearson's chi-square test statistic. We see that the LRT statistic $2 \sum_{i=1}^k Y_i \log \frac{\hat{p}_i}{p_i^*}$ and Pearson's statistic asymptotically are equivalent under H_0 . Pearson's statistic can be viewed as quadratic approximation of the LRT statistic.

We can write Pearson's statistic as

$$\sum_{i=1}^k \frac{(Y_i - Np_i^*)^2}{Np_i^*} = N(\hat{\mathbf{p}} - \mathbf{p}^*)' \mathbf{Q}(\hat{\mathbf{p}} - \mathbf{p}^*),$$

where $\hat{\mathbf{p}} = (Y_1/N, \dots, Y_k/N)'$ and $\mathbf{Q} := \text{diag}(1/p_1^*, \dots, 1/p_k^*)$. By the CLT, under H_0 , $\sqrt{N}(\hat{\mathbf{p}} - \mathbf{p}^*)$ converges in distribution to normal $\mathcal{N}_k(\mathbf{0}, \mathbf{C})$. Recall that since $\mathbf{1}_k' \hat{\mathbf{p}} = 1$ and $\mathbf{1}_k' \mathbf{p}^* = 1$, the

covariance matrix \mathbf{C} has rank $k - 1$, and hence is singular. Therefore the normal distribution $\mathcal{N}_k(\mathbf{0}, \mathbf{C})$ is degenerate.

Consider $\mathbf{Z} \sim \mathcal{N}_k(\mathbf{0}, \mathbf{C})$, let us show that $\mathbf{Z}'\mathbf{Q}\mathbf{Z}$ has χ_{k-1}^2 distribution. For $\mathbf{W} := \mathbf{Q}^{1/2}\mathbf{Z}$ we have that $\mathbf{W}'\mathbf{W} = \mathbf{Z}'\mathbf{Q}\mathbf{Z}$ and $\mathbf{W} \sim \mathcal{N}_k(\mathbf{0}, \mathbf{M})$, where $\mathbf{M} := \mathbf{I}_k - (\mathbf{p}^{*1/2})(\mathbf{p}^{*1/2})'$. Matrix \mathbf{M} is a projection matrix of rank

$$\text{rank}(\mathbf{M}) = \text{tr}(\mathbf{I}_k - (\mathbf{p}^{*1/2})(\mathbf{p}^{*1/2})') = k - (\mathbf{p}^{*1/2})'(\mathbf{p}^{*1/2}) = k - \sum_{i=1}^k p_i^* = k - 1.$$

Since \mathbf{M} is a projection matrix of rank $k - 1$, it has $k - 1$ eigenvalues equal 1 and one eigenvalue 0. Therefore it has the spectral decomposition $\mathbf{M} = \mathbf{T}\mathbf{\Lambda}\mathbf{T}'$ with $\mathbf{\Lambda} = \text{diag}(1, \dots, 1, 0)$. Consider $\mathbf{Y} := \mathbf{T}'\mathbf{W}$. Since matrix \mathbf{T} is orthogonal we have that $\mathbf{W}'\mathbf{W} = \mathbf{Y}'\mathbf{Y}$. Also $\mathbf{Y} \sim \mathcal{N}_k(\mathbf{0}, \mathbf{T}'\mathbf{M}\mathbf{T})$. Since the last element of matrix $\mathbf{T}'\mathbf{M}\mathbf{T} = \mathbf{\Lambda}$ is zero it follows that $\text{Var}(Y_k) = 0$ and hence $Y_k \equiv 0$. Therefore $\mathbf{Y}'\mathbf{Y} = Y_1^2 + \dots + Y_{k-1}^2 \sim \chi_{k-1}^2$. It follows that under H_0 , $N(\hat{\mathbf{p}} - \mathbf{p}^*)'\mathbf{Q}(\hat{\mathbf{p}} - \mathbf{p}^*)$ converges in distribution to χ_{k-1}^2 . \square

General model: $\mathbf{p} = \mathbf{p}(\boldsymbol{\theta})$, $\boldsymbol{\theta} \in \mathbb{R}^q$, with $\sum_{i=1}^k p_i(\boldsymbol{\theta}) = 1$. The ML estimator of parameter vector $\boldsymbol{\theta}$ is solution of the optimization problem

$$\max_{\boldsymbol{\theta}} \sum_{i=1}^k Y_i \log p_i(\boldsymbol{\theta}).$$

Suppose that the model is correct, i.e., there is $\boldsymbol{\theta}^* \in \mathbb{R}^q$ such that $\mathbf{p}^* = \mathbf{p}(\boldsymbol{\theta}^*)$, where \mathbf{p}^* is the true (population) value of the parameter vector. Suppose further that the model is identified at $\boldsymbol{\theta}^*$, i.e., if $\mathbf{p}(\boldsymbol{\theta}) = \mathbf{p}(\boldsymbol{\theta}^*)$, then $\boldsymbol{\theta} = \boldsymbol{\theta}^*$. Let $\tilde{\mathbf{p}}$ be a consistent estimator of \mathbf{p}^* . Then asymptotically the ML estimation is equivalent to

$$\min_{\boldsymbol{\theta}} (\tilde{\mathbf{p}} - \mathbf{p}(\boldsymbol{\theta}))'\tilde{\mathbf{Q}}(\tilde{\mathbf{p}} - \mathbf{p}(\boldsymbol{\theta})),$$

where $\hat{p}_i = Y_i/N$, $i = 1, \dots, N$ and $\tilde{\mathbf{Q}} = \text{diag}(1/\tilde{p}_1, \dots, 1/\tilde{p}_k)$.

We have here that $\sqrt{N}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)$ converges in distribution to normal $\mathcal{N}(\mathbf{0}, \mathbf{I}(\boldsymbol{\theta}^*)^{-1})$ with $\mathbf{I}(\boldsymbol{\theta}) = \mathbf{P}(\boldsymbol{\theta})'\mathbf{C}(\boldsymbol{\theta})\mathbf{P}(\boldsymbol{\theta})$, where $\mathbf{P}(\boldsymbol{\theta}) = \partial \log \mathbf{p}(\boldsymbol{\theta})/\partial \boldsymbol{\theta}'$ is $m \times q$ matrix and $\mathbf{C}(\boldsymbol{\theta}) = \mathbf{P}(\boldsymbol{\theta}) - \mathbf{p}(\boldsymbol{\theta})\mathbf{p}(\boldsymbol{\theta})'$. The LRT for testing $H_0 : \mathbf{p} = \mathbf{p}(\boldsymbol{\theta})$ against the unrestricted alternative is

$$2 \log \lambda = 2 \sum_{i=1}^k Y_i \log \frac{Y_i/N}{p_i(\tilde{\boldsymbol{\theta}})},$$

where $\tilde{\boldsymbol{\theta}}$ is the MLE under H_0 . Under H_0 , the LRT statistic $2 \log \lambda$ converges in distribution to χ_{k-1-q}^2 , and asymptotically is equivalent to Pearson's statistic.

11 Logistic regression

Let Y_1, \dots, Y_N be independent random variables such that Y_i has the binomial distribution $B(m_i, \pi_i)$, $i = 1, \dots, N$. Consider the logit model:

$$\pi_i = \frac{\exp(\beta_0 + \beta_1 X_{i1} + \dots + \beta_k X_{ik})}{1 + \exp(\beta_0 + \beta_1 X_{i1} + \dots + \beta_k X_{ik})}, \quad i = 1, \dots, N, \quad (11.1)$$

where β_0, \dots, β_k are parameters. That is

$$\log \frac{\pi_i}{1 - \pi_i} = \beta_0 + \beta_1 X_{i1} + \dots + \beta_k X_{ik}, \quad i = 1, \dots, N,$$

where $\frac{\pi_i}{1-\pi_i}$ is called the *odds ratio*.

We can write this model in the matrix form

$$\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}, \quad (11.2)$$

where $\eta_i := \log \frac{\pi_i}{1-\pi_i}$, $i = 1, \dots, N$, and $\mathbf{X} = [\mathbf{1}_N, \mathbf{X}_1, \dots, \mathbf{X}_k]$ is the design matrix. As in the linear regression we assume that matrix \mathbf{X} has full column rank $p = k + 1$. The multicollinearity problem can also happen here when columns of the design matrix are ‘almost’ linearly dependent.

We have that

$$P(Y_i = y_i) = \binom{m_i}{y_i} \pi_i^{y_i} (1 - \pi_i)^{m_i - y_i}, \quad y_i = 0, 1, \dots, m_i.$$

It follows that the likelihood function here is

$$L(\boldsymbol{\pi}; \mathbf{y}) = c \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{m_i - y_i},$$

where the constant $c = \prod_{i=1}^n \binom{m_i}{y_i}$ is independent of $\boldsymbol{\pi}$. Hence up to the constant $\log c$ independent of $\boldsymbol{\pi}$, the log likelihood function $\log L(\boldsymbol{\pi}; \mathbf{y})$ can be written as

$$l(\boldsymbol{\pi}; \mathbf{y}) = \sum_{i=1}^n [y_i \log \pi_i + (m_i - y_i) \log(1 - \pi_i)],$$

(note that, by definition, $0 \times \log 0 = 0$).

Fisher’s information matrix, for $\boldsymbol{\beta} = (\beta_0, \dots, \beta_k)'$, can be written in the form $\mathbf{X}'\mathbf{W}\mathbf{X}$, where \mathbf{W} is a diagonal matrix given by

$$\mathbf{W} = \text{diag}\{m_1\pi_1(1 - \pi_1), \dots, m_n\pi_n(1 - \pi_n)\}.$$

Indeed, we have that

$$\frac{\partial l}{\partial \pi_i} = \frac{y_i - m_i\pi_i}{\pi_i(1 - \pi_i)},$$

and hence

$$\frac{\partial l}{\partial \beta_s} = \sum_{i=1}^n \frac{y_i - m_i\pi_i}{\pi_i(1 - \pi_i)} \frac{\partial \pi_i}{\partial \beta_s},$$

where

$$\frac{\partial \pi_i}{\partial \beta_s} = \pi_i(1 - \pi_i)X_{si}.$$

Consequently the st -element of Fisher’s information matrix is

$$\mathbb{E}\left[\frac{\partial l}{\partial \beta_s} \frac{\partial l}{\partial \beta_t}\right] = \mathbb{E}\left[\sum_{i,j} \left(\frac{Y_i - m_i\pi_i}{\pi_i(1 - \pi_i)} \frac{\partial \pi_i}{\partial \beta_s}\right) \left(\frac{Y_j - m_j\pi_j}{\pi_j(1 - \pi_j)} \frac{\partial \pi_j}{\partial \beta_t}\right)\right], \quad s, t = 0, \dots, k.$$

Moreover, $\mathbb{E}[Y_i] = m_i\pi_i$, and hence (by independence)

$$\mathbb{E}[(Y_i - m_i\pi_i)(Y_j - m_j\pi_j)] = 0, \quad \text{if } i \neq j,$$

and

$$\mathbb{E}[(Y_i - m_i\pi_i)^2] = \text{Var}[Y_i] = m_i\pi_i(1 - \pi_i), \quad i = 1, \dots, n.$$

It follows that

$$\mathbb{E} \left[\frac{\partial l}{\partial \beta_s} \frac{\partial l}{\partial \beta_t} \right] = \sum_{i=1}^n \frac{m_i}{\pi_i(1-\pi_i)} \frac{\partial \pi_i}{\partial \beta_s} \frac{\partial \pi_i}{\partial \beta_t} = \sum_{i=1}^n m_i \pi_i (1-\pi_i) X_{si} X_{ti}.$$

The maximum likelihood (ML) equations are

$$\sum_{i=1}^n (y_i - m_i \pi_i) X_{si} = 0, \quad s = 0, \dots, k.$$

Consider the log-likelihood function $l(\cdot; \mathbf{y})$ as a function of \mathbf{x} with $\boldsymbol{\pi} = \boldsymbol{\pi}(\mathbf{x})$. We have that the terms $y_i \log \pi_i$ are linear functions of \mathbf{x} , the terms $(m_i - y_i) \log(1 - \pi_i)$ consist of linear terms and terms of the form $-(m_i - y_i) \log(1 + \exp(1 + \boldsymbol{\beta}'\mathbf{x}))$. Since the function $\phi(x) = \log(1 + e^x)$ is strictly convex, it follows that $l(\cdot; \mathbf{y})$ is strictly concave function of \mathbf{x} , and hence the ML equations for estimating $\boldsymbol{\beta}$ have *unique* solution $\hat{\boldsymbol{\beta}}$ (recall that the design matrix in (11.2) is assumed to have full column rank).

Consider

$$\hat{\pi}_i = \frac{\exp(\hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \dots + \hat{\beta}_k X_{ik})}{1 + \exp(\hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \dots + \hat{\beta}_k X_{ik})}, \quad i = 1, \dots, n,$$

and the following so-called deviance function,

$$\Lambda = -2l(\tilde{\boldsymbol{\pi}}; \mathbf{y}) + 2l(\hat{\boldsymbol{\pi}}; \mathbf{y}),$$

where $\tilde{\boldsymbol{\pi}}$ is the ML estimate of $\boldsymbol{\pi}$ under a specified H_0 . That is, Λ is the log-likelihood ratio test statistic $2 \log \lambda$ for testing H_0 . In particular, for $H_0 : \beta_1 = \dots = \beta_k = 0$ we have that $\tilde{\pi}_i = \tilde{\pi}$, $i = 1, \dots, N$, where $\tilde{\pi} = \frac{\sum_{i=1}^N y_i}{\sum_{i=1}^N m_i}$.

If $m_i = 1$, $i = 1, \dots, N$, then Y_1, \dots, Y_N become Bernoulli random variables with $P(Y_i = 1) = \pi_i$ and $P(Y_i = 0) = 1 - \pi_i$. In that case

$$l(\boldsymbol{\pi}; \mathbf{y}) = \sum_{i=1}^N [y_i \log \pi_i + (1 - y_i) \log(1 - \pi_i)].$$

For $H_0 : \beta_1 = \dots = \beta_k = 0$ we have that $\tilde{\pi}_i = \tilde{\pi}$, $i = 1, \dots, N$, where $\tilde{\pi} = \frac{\sum_{i=1}^N y_i}{N}$, and hence $l(\boldsymbol{\pi}; \mathbf{y}) = \left(\sum_{i=1}^N y_i \right) \log \tilde{\pi}$.

12 Generalized linear models

Let $\mathbf{Y} = (Y_1, \dots, Y_N)'$ be a vector of responses whose components are independently distributed with means $\boldsymbol{\mu} = (\mu_1, \dots, \mu_N)'$, i.e., $\mu_i = \mathbb{E}[Y_i]$, $i = 1, \dots, N$. The linear model assumes that $\mu_i = \mathbf{x}'_i \boldsymbol{\beta}$, $i = 1, \dots, N$, where $\boldsymbol{\beta}$ is $k \times 1$ vector of parameters and $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ik})'$ are observed values of the predictors. That is, the conditional expectation $\mathbb{E}[Y_i | \mathbf{X}_i = \mathbf{x}_i] = \mathbf{x}'_i \boldsymbol{\beta}$, $i = 1, \dots, N$.

This can be generalized in the following way. Let us introduce a linear predictor

$$\eta_i = \mathbf{x}'_i \boldsymbol{\beta}, \quad i = 1, \dots, N. \quad (12.1)$$

The new symbol η is related to μ by the equation $\eta = g(\mu)$, where $g(\cdot)$ is a specified function called the *link function*. That is

$$\eta_i = g(\mu_i), \quad i = 1, \dots, N,$$

and

$$\mu_i = g^{-1}(\mathbf{x}'_i \boldsymbol{\beta}), \quad i = 1, \dots, N.$$

We also can write it in the matrix form

$$\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta},$$

where \mathbf{X} is the design matrix with rows $(1, X_{i1}, \dots, X_{ik})$, $i = 1, \dots, N$. As before it is assumed that the design matrix has full column rank.

For example, in the linear case $\boldsymbol{\eta} = \boldsymbol{\mu}$, i.e., $g(\mu) = \mu$. In the logistic regression $g(\pi) = \log \frac{\pi}{1-\pi}$ is the logit link function.

Suppose now that each component Y_i of the response vector has a distribution in the exponential family with pdf of the form

$$f_Y(y; \theta, \phi) = \exp \left\{ \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right\} \quad (12.2)$$

for some specified functions $a(\cdot)$, $b(\cdot)$ and $c(\cdot)$. The parameter θ is called the natural parameter, and the parameter ϕ the dispersion parameter. For example for the normal distribution $\mathcal{N}(\mu, \sigma^2)$ we can write the corresponding density

$$f(y) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{(y - \mu)^2}{2\sigma^2} \right)$$

in the form (12.2) with $\theta = \mu$, $\phi = \sigma$ and

$$a(\phi) = \phi^2, \quad b(\theta) = \theta^2/2, \quad c(y, \phi) = -\frac{1}{2}\{y^2/\phi^2 + \log(2\pi\phi^2)\}.$$

If ϕ is known, then $a(\phi)$ is viewed as a constant, $c(y, \phi) = c(y)$, and (12.2) becomes an exponential family in the canonical form with canonical parameter θ .

Consider

$$l(y; \theta, \phi) = \log f_Y(y, \theta, \phi) = \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi). \quad (12.3)$$

By the standard theory of the ML we have that

$$\mathbb{E}[\partial l / \partial \theta] = 0, \quad (12.4)$$

$$\mathbb{E}[\partial^2 l / \partial \theta^2] = -\mathbb{E}[(\partial l / \partial \theta)^2]. \quad (12.5)$$

Also by (12.3)

$$\partial l / \partial \theta = \frac{y - b'(\theta)}{a(\phi)}$$

and because of (12.4), $\mathbb{E}[Y - b'(\theta)] = 0$. Thus $\mathbb{E}[Y] = b'(\theta)$, that is (compare with (7.2))

$$\mu = b'(\theta).$$

Moreover

$$\partial^2 l / \partial \theta^2 = -b''(\theta) / a(\phi)$$

and hence $b''(\theta) / a(\phi) = \text{Var}(Y) / a^2(\phi)$ and thus (compare with (7.3))

$$\text{Var}(Y) = b''(\theta) a(\phi).$$

For binomial distribution $B(m, \pi)/m$ the corresponding distribution function is

$$P(Y = y) = \binom{m}{my} \pi^{my} (1 - \pi)^{m(1-y)}, \quad y = 0, 1/m, \dots, 1.$$

Let us set $\theta = \log \frac{\pi}{1-\pi}$ as the natural parameter, and hence $\pi = \frac{e^\theta}{1+e^\theta}$. Here $\mu = \pi$ and thus $\mu = \frac{e^\theta}{1+e^\theta}$. Assume that m is known and set $\phi = 1/m$, $a(\phi) = \phi$, $b(\theta) = \log(1 + e^\theta)$, $c(y, \phi) = \log \binom{m}{my}$. Note that $0 \log 0 = 0$, and hence for $m = 1$ we have that $\phi = 1$ and $c(y, \phi) = 0$. The link function here is $\text{logit } g(\pi) = \log \frac{\pi}{1-\pi}$.

For Poisson distribution

$$P(Y = y) = \frac{1}{y!} e^{-\mu} \mu^y, \quad y = 0, 1, 2, \dots,$$

with parameter $\mu > 0$. Note that $\mu = \mathbb{E}[Y]$ here. This can be written as

$$P(Y = y) = \exp\{y \log \mu - \mu - \log(y!)\}, \quad y = 0, 1, 2, \dots$$

We have here that $\mu = \mathbb{E}[Y]$ and $\theta = \log \mu$ is the natural parameter with $b(\theta) = e^\theta$, $a(\phi) = 1$ and $c(y) = -\log(y!)$. The link function here is $g(\mu) = \log \mu$.

In the canonical case (when ϕ is known) the model is $\theta_i = \eta_i$, $i = 1, \dots, n$, with η_i being linear predictors specified in equation (12.1). In order to compute the ML estimate of β we need to maximize the corresponding log-likelihood function (given in (12.3)), that is to solve the problem

$$\max_{\beta} \sum_{i=1}^n Y_i \mathbf{x}'_i \beta - b(\mathbf{x}'_i \beta). \quad (12.6)$$

When $b(\cdot)$ is a convex function, the above problem (12.6) is convex. For the binomial and Poisson distributions the corresponding functions $b(\cdot)$ are convex.

13 Classification problem

Consider an $m \times 1$ random vector \mathbf{X} of measurements. We want to classify \mathbf{X} into one of two population π_1 or π_2 . Let $p_1(\mathbf{x})$ and $p_2(\mathbf{x})$ be respective densities (pdfs) of populations π_1 and π_2 . Suppose that the probability that an observation comes from population π_i is q_i , $i = 1, 2$. Consider regions $R_1 \subset \mathbb{R}^m$ and $R_2 = \mathbb{R}^m \setminus R_1$. If $\mathbf{X} \in R_1$ we classify \mathbf{X} as from π_1 , and if $\mathbf{X} \in R_2$ we classify \mathbf{X} as from π_2 . Then the probability of misclassification of an observation from π_1 is

$$\text{Prob}(\mathbf{X} \in R_2 | \pi_1) = \int_{R_2} p_1(\mathbf{x}) d\mathbf{x}.$$

Similarly the probability of misclassification of an observation from π_2 is $\int_{R_1} p_2(\mathbf{x}) d\mathbf{x}$. The expected loss of misclassification is

$$c_1 q_1 \int_{R_2} p_1(\mathbf{x}) d\mathbf{x} + c_2 q_2 \int_{R_1} p_2(\mathbf{x}) d\mathbf{x},$$

where c_i is the cost of misclassification of an observation from π_i , $i = 1, 2$.

Note that

$$\int_{R_1} p_2(\mathbf{x}) d\mathbf{x} = \int_{\mathbb{R}^m \setminus R_2} p_2(\mathbf{x}) d\mathbf{x} = \int_{\mathbb{R}^m} p_2(\mathbf{x}) d\mathbf{x} - \int_{R_2} p_2(\mathbf{x}) d\mathbf{x}.$$

Suppose that $c_1 = c_2 = 1$. Then the probability (expected loss) of misclassification is

$$q_1 \int_{R_2} p_1(\mathbf{x}) d\mathbf{x} + q_2 \int_{R_1} p_2(\mathbf{x}) d\mathbf{x} = \int_{R_2} [q_1 p_1(\mathbf{x}) - q_2 p_2(\mathbf{x})] d\mathbf{x} + q_2 \int_{\mathbb{R}^m} p_2(\mathbf{x}) d\mathbf{x}.$$

Since $p_2(\cdot)$ is a probability density function, we have that $\int_{\mathbb{R}^m} p_2(\mathbf{x}) d\mathbf{x} = 1$, and hence

$$q_1 \int_{R_2} p_1(\mathbf{x}) d\mathbf{x} + q_2 \int_{R_1} p_2(\mathbf{x}) d\mathbf{x} = \int_{R_2} [q_1 p_1(\mathbf{x}) - q_2 p_2(\mathbf{x})] d\mathbf{x} + q_2.$$

It follows that the expected loss is minimized if

$$R_2 = \{\mathbf{x} \in \mathbb{R}^m : q_1 p_1(\mathbf{x}) - q_2 p_2(\mathbf{x}) < 0\}.$$

Or equivalently

$$R_1 = \left\{ \mathbf{x} \in \mathbb{R}^m : p_1(\mathbf{x}) \geq \frac{q_2}{q_1} p_2(\mathbf{x}) \right\}$$

and

$$R_2 = \left\{ \mathbf{x} \in \mathbb{R}^m : p_1(\mathbf{x}) < \frac{q_2}{q_1} p_2(\mathbf{x}) \right\}.$$

If the costs c_1 and c_2 are unequal, then the optimal regions are

$$R_1 = \left\{ \mathbf{x} \in \mathbb{R}^m : p_1(\mathbf{x}) \geq \frac{c_2 q_2}{c_1 q_1} p_2(\mathbf{x}) \right\}$$

and

$$R_2 = \left\{ \mathbf{x} \in \mathbb{R}^m : p_1(\mathbf{x}) < \frac{c_2 q_2}{c_1 q_1} p_2(\mathbf{x}) \right\}.$$

If $c_1 q_1 p_1(\mathbf{x}) = c_2 q_2 p_2(\mathbf{x})$, we can take such points either in R_1 or R_2 .

It could be noted that the above derivations basically are the same as derivation of the Neyman - Pearson Lemma in section 9. The only difference is that the misclassification errors are treated here symmetrically, unlike in the hypothesis testing.

13.1 Classification with normally distributed populations

Suppose that the populations π_1 and π_2 have multivariate normal distributions with equal covariance matrices, i.e., $\pi_i \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$, $i = 1, 2$. Then

$$p_i(\mathbf{x}) = \frac{1}{(2\pi)^{m/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -(\mathbf{x} - \boldsymbol{\mu}_i)' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) / 2 \right\},$$

and

$$\frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} = \exp \left\{ -\frac{1}{2} [(\mathbf{x} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) - (\mathbf{x} - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)] \right\}.$$

Hence the optimal region is

$$R_1 = \left\{ \mathbf{x} : (\mathbf{x} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) - (\mathbf{x} - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \leq -2\kappa \right\},$$

where $\kappa = \log(c_2 q_2 / c_1 q_1)$. Equivalently

$$R_1 = \left\{ \mathbf{x} : \mathbf{x}' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \geq \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + \kappa \right\}. \quad (13.1)$$

Note that if $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$, then $\mathbf{X}' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ has normal distribution with mean $\boldsymbol{\mu}_i' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ and variance $(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$. The function $\mathbf{X}' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ is called *Fisher's discriminant function*, and $\sqrt{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)}$ is called *Mahalanobis' distance* between $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$.

13.1.1 An optimization problem

Before proceeding further we need the following result. Consider optimization problem

$$\max_{\mathbf{d} \neq \mathbf{0}} \frac{\mathbf{d}' \mathbf{A} \mathbf{d}}{\mathbf{d}' \mathbf{B} \mathbf{d}}, \quad (13.2)$$

where \mathbf{A} is an $m \times m$ symmetric positive semidefinite matrix and \mathbf{B} is an $m \times m$ symmetric positive definite matrix. Let $\mathbf{B}^{1/2}$ be symmetric positive definite matrix such that $\mathbf{B} = \mathbf{B}^{1/2} \mathbf{B}^{1/2}$ (see section 1 for discussion of such functions of symmetric matrices). By change of variables $\mathbf{h} = \mathbf{B}^{1/2} \mathbf{d}$ we can write problem (13.2) as

$$\max_{\mathbf{h} \neq \mathbf{0}} \frac{\mathbf{h}' (\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2}) \mathbf{h}}{\|\mathbf{h}\|^2},$$

where $\|\cdot\|$ denotes the Euclidean norm. This in turn can be written as

$$\max_{\|\mathbf{h}\|=1} \mathbf{h}' (\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2}) \mathbf{h}. \quad (13.3)$$

Such problems are discussed in section 15.

Matrix $\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2}$ is symmetric positive semidefinite. Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ be the eigenvalues and $\mathbf{e}_1, \dots, \mathbf{e}_m$ be the corresponding orthonormal eigenvectors of matrix $\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2}$, i.e., $\|\mathbf{e}_i\| = 1$ and $\mathbf{e}_i' \mathbf{e}_j = 0$ for $i \neq j$. Note that $\mathbf{B}^{-1/2} \mathbf{e}_i$ are eigenvectors of matrix $\mathbf{B}^{-1} \mathbf{A}$ corresponding to the same eigenvalues λ_i . This follows from $\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2} \mathbf{e}_i = \lambda_i \mathbf{e}_i$ by multiplying both sides of this equation by $\mathbf{B}^{-1/2}$. The optimal solution $\bar{\mathbf{h}}$ of problem (13.3) is the eigenvector $\bar{\mathbf{h}} = \mathbf{e}_1$. It follows that solution $\mathbf{d}_1 = \mathbf{B}^{-1/2} \mathbf{e}_1$ of problem (13.2) is given by the eigenvector of $\mathbf{B}^{-1} \mathbf{A}$ corresponding to its largest eigenvalue. Note that the optimal solution of problem (13.2) is defined up to a scale change, i.e., changing \mathbf{d} to $t\mathbf{d}$ does not change value of the objective function in (13.2) for any nonzero number t .

In particular, suppose that matrix \mathbf{A} has rank one, and hence can be written as $\mathbf{A} = \mathbf{a} \mathbf{a}'$ where $\mathbf{a} \neq \mathbf{0}$ is an $m \times 1$ vector. Then $\mathbf{d}_1 = \mathbf{B}^{-1} \mathbf{a}$ is an optimal solution of problem (13.2). Indeed, in that case matrix $\mathbf{B}^{-1} \mathbf{A}$ has only one nonzero eigenvalue, the largest one. Moreover, for $\lambda = \mathbf{a}' \mathbf{B}^{-1} \mathbf{a}$ we have

$$\mathbf{B}^{-1} \mathbf{A} \mathbf{d}_1 = \mathbf{B}^{-1} \mathbf{a} (\mathbf{a}' \mathbf{B}^{-1} \mathbf{a}) = \lambda \mathbf{B}^{-1} \mathbf{a} = \lambda \mathbf{d}_1. \quad (13.4)$$

It also follows that in that case the optimal value of problem (13.2) is equal to $\mathbf{a}' \mathbf{B}^{-1} \mathbf{a}$.

Next consider the following problem

$$\max_{\mathbf{d} \neq \mathbf{0}} \frac{\mathbf{d}' \mathbf{A} \mathbf{d}}{\mathbf{d}' \mathbf{B} \mathbf{d}} \text{ subject to } \mathbf{d}' \mathbf{B} \mathbf{d}_1 = 0. \quad (13.5)$$

Again by change of variables $\mathbf{h} = \mathbf{B}^{1/2} \mathbf{d}$ we obtain the problem

$$\max_{\|\mathbf{h}\|=1} \mathbf{h}' (\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2}) \mathbf{h}, \text{ subject to } \mathbf{h}' \mathbf{e}_1 = 0.$$

Optimal solution of this problem is the eigenvector \mathbf{e}_2 of $\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2}$, and hence the optimal solution of problem (13.5) is the eigenvector $\mathbf{d}_2 = \mathbf{B}^{-1/2} \mathbf{e}_2$ of $\mathbf{B}^{-1} \mathbf{A}$ corresponding to its second largest eigenvalue. Note that $\mathbf{d}_2' \mathbf{B} \mathbf{d}_1 = \mathbf{e}_2' \mathbf{e}_1 = 0$.

13.2 Fisher discriminant analysis

Suppose that distribution of population π_i has mean $\boldsymbol{\mu}_i$ and covariance matrix $\boldsymbol{\Sigma}_i$. Consider the following problem

$$\max_{\mathbf{d} \neq \mathbf{0}} \left\{ g(\mathbf{d}) := \frac{(\mathbf{d}'\boldsymbol{\mu}_1 - \mathbf{d}'\boldsymbol{\mu}_2)^2}{\mathbf{d}'\boldsymbol{\Sigma}_1\mathbf{d} + \mathbf{d}'\boldsymbol{\Sigma}_2\mathbf{d}} \right\}. \quad (13.6)$$

Note that $\mathbf{d}'\boldsymbol{\mu}_i$ is the expected value and $\mathbf{d}'\boldsymbol{\Sigma}_i\mathbf{d}$ is the variance of $\mathbf{d}'\mathbf{X}$ for population π_i .

We can write function $g(\mathbf{d})$ as

$$g(\mathbf{d}) = \frac{\mathbf{d}'(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)'\mathbf{d}}{\mathbf{d}'(\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)\mathbf{d}}.$$

Hence the optimal solution $\bar{\mathbf{d}}$ of problem (13.6) is (see equation (13.4))

$$\bar{\mathbf{d}} = (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2). \quad (13.7)$$

In particular if $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}$, then $\bar{\mathbf{d}} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$. Recall that the optimal solution of problem (13.6) is defined up to a scale change.

13.3 Several populations

Suppose that there are r populations π_1, \dots, π_r with respective means $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_r$ and covariance matrices $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_r$. Let q_i be the probability that the measurements vector \mathbf{X} comes from population π_i , $i = 1, \dots, r$ (we assume that $q_i > 0$, $i = 1, \dots, r$). We have that

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}] = q_1\boldsymbol{\mu}_1 + \dots + q_r\boldsymbol{\mu}_r$$

and

$$\mathbb{E}[\mathbf{X}\mathbf{X}'] = q_1(\boldsymbol{\Sigma}_1 + \boldsymbol{\mu}_1\boldsymbol{\mu}_1') + \dots + q_r(\boldsymbol{\Sigma}_r + \boldsymbol{\mu}_r\boldsymbol{\mu}_r'),$$

and hence

$$\text{Cov}(\mathbf{X}) = \mathbb{E}[\mathbf{X}\mathbf{X}'] - \boldsymbol{\mu}\boldsymbol{\mu}' = \sum_{i=1}^r q_i\boldsymbol{\Sigma}_i + \sum_{i=1}^r q_i(\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})' = \boldsymbol{\Omega} + \mathbf{M},$$

where

$$\boldsymbol{\Omega} := \sum_{i=1}^r q_i\boldsymbol{\Sigma}_i \quad \text{and} \quad \mathbf{M} := \sum_{i=1}^r q_i(\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})'. \quad (13.8)$$

Consider the following optimization problem

$$\max_{\mathbf{d} \in \mathbb{R}^m} \left\{ g(\mathbf{d}) := \frac{\mathbf{d}'\mathbf{M}\mathbf{d}}{\mathbf{d}'\boldsymbol{\Omega}\mathbf{d}} \right\}. \quad (13.9)$$

Note that matrices $\boldsymbol{\Sigma}_i$ are positive definite and hence matrix $\boldsymbol{\Omega}$ is positive definite, and matrix \mathbf{M} is positive semidefinite. Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ be the eigenvalues of $\boldsymbol{\Omega}^{-1}\mathbf{M}$. Then the optimal solution \mathbf{d}_1 of problem (13.9) is the eigenvector of $\boldsymbol{\Omega}^{-1}\mathbf{M}$ corresponding to its largest eigenvalue λ_1 (see section 13.1.1). Next maximize $g(\mathbf{d})$ subject to $\mathbf{d}'\boldsymbol{\Omega}\mathbf{d}_1 = 0$. The solution of this problem is given by eigenvector \mathbf{d}_2 of $\boldsymbol{\Omega}^{-1}\mathbf{M}$ corresponding to the second largest eigenvalue λ_2 . By continuing this process we obtain discriminant functions $\mathbf{d}'_i\mathbf{X}$, $i = 1, \dots, r-1$. Note that $\text{rank}(\mathbf{M}) \leq r-1$ since

$$\sum_{i=1}^r q_i(\boldsymbol{\mu}_i - \boldsymbol{\mu}) = \sum_{i=1}^r q_i\boldsymbol{\mu}_i - \boldsymbol{\mu} = \mathbf{0}.$$

Hence $\lambda_r = \dots = \lambda_m = 0$. For $r = 2$ we have $\boldsymbol{\mu}_1 - \boldsymbol{\mu} = q_2(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$, $\boldsymbol{\mu}_2 - \boldsymbol{\mu} = q_1(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ and hence $\mathbf{M} = q_1q_2(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)'$. In that case the above approach is the same as Fisher's discriminant analysis.

13.3.1 Mahalanobis distance

Mahalanobis distance between two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^m$, with respect to covariance matrix Σ , is defined as

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})' \Sigma^{-1} (\mathbf{x} - \mathbf{y})}$$

Assuming that covariance matrices $\Sigma_1 = \dots = \Sigma_r = \Sigma$ are equal to each other, classify \mathbf{X} in π_i if $d(\mathbf{X}, \boldsymbol{\mu}_i) < d(\mathbf{X}, \boldsymbol{\mu}_j)$ for all $j \neq i$.

Voronoi diagram. The positive definite matrix Σ^{-1} defines the corresponding norm $\|\mathbf{x}\|_{\Sigma^{-1}} := \sqrt{\mathbf{x}' \Sigma^{-1} \mathbf{x}}$. If $\Sigma = \mathbf{I}_m$ this is the Euclidean norm.

Partition of \mathbb{R}^m into regions

$$R_i = \{\mathbf{x} : \|\mathbf{x} - \boldsymbol{\mu}_i\|_{\Sigma^{-1}} \leq \|\mathbf{x} - \boldsymbol{\mu}_j\|_{\Sigma^{-1}}, j \neq i\}, \quad i = 1, \dots, r,$$

is called Voronoi diagram (with respect to the norm $\|\cdot\|_{\Sigma^{-1}}$). Note that each set R_i is polyhedral given by intersection of half spaces

$$\{\mathbf{x} : \mathbf{x}' \Sigma^{-1} (\boldsymbol{\mu}_j - \boldsymbol{\mu}_i) \leq \frac{1}{2} (\boldsymbol{\mu}_j' \Sigma^{-1} \boldsymbol{\mu}_j - \boldsymbol{\mu}_i' \Sigma^{-1} \boldsymbol{\mu}_i)\}, \quad j \neq i.$$

Mahalanobis distance classification: classify \mathbf{X} in π_i if $\mathbf{X} \in R_i$. For $r = 2$ this is the same classification as in (13.1) with $q_1 = q_2$ and $c_1 = c_2$.

13.4 Bayes and Logistic Regression classifiers

Suppose that we have two populations π_1 and π_2 . We consider (Y, \mathbf{X}) with $Y = 1$ if $\mathbf{X} \sim \pi_1$ and $Y = -1$ if $\mathbf{X} \sim \pi_2$. By Bayes formula we have that

$$\text{Prob}(Y = 1 | \mathbf{X} = \mathbf{x}) = \frac{p_1(\mathbf{x})q_1}{p_1(\mathbf{x})q_1 + p_2(\mathbf{x})q_2},$$

where $q_1 = \text{Prob}(Y = 1)$ and $q_2 = \text{Prob}(Y = -1)$. We classify \mathbf{X} in π_1 if $\text{Prob}(Y = 1 | \mathbf{X} = \mathbf{x}) > \text{Prob}(Y = -1 | \mathbf{X} = \mathbf{x})$, which is equivalent to $p_1(\mathbf{x})q_1 > p_2(\mathbf{x})q_2$.

Logistic regression approach. The ratio $\text{odd}(\mathbf{x}) = \frac{\text{Prob}(Y=1|\mathbf{X}=\mathbf{x})}{\text{Prob}(Y=-1|\mathbf{X}=\mathbf{x})}$ is called odds ratio. Logistic regression model:

$$\log \text{odd}(\mathbf{x}) = \beta_0 + \boldsymbol{\beta}' \mathbf{x}. \quad (13.10)$$

We classify \mathbf{X} in π_1 if $\text{odd}(\mathbf{x}) > 1$. This is equivalent to $\beta_0 + \boldsymbol{\beta}' \mathbf{x} > 0$.

Note that

$$\frac{\text{Prob}(Y = 1 | \mathbf{X} = \mathbf{x})}{\text{Prob}(Y = -1 | \mathbf{X} = \mathbf{x})} = \frac{p_1(\mathbf{x})q_1}{p_2(\mathbf{x})q_2}.$$

In case of normal distributions with the same covariance matrix Σ we have that (see section 13.1)

$$\frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} = \exp \{ \mathbf{x}' \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + \text{const} \}.$$

In that case (assuming $q_1 = q_2$) equation (13.10) holds with $\boldsymbol{\beta} = \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$.

14 Support Vector Machines

Suppose that we have two populations π_1 and π_2 . Suppose further that we have training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$, where $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i = 1$ if $\mathbf{x}_i \sim \pi_1$ and $y_i = -1$ if $\mathbf{x}_i \sim \pi_2$. We want to separate these populations by a hyperplane $\beta_0 + \boldsymbol{\beta}'\mathbf{x} = 0$. That is, we classify an observation \mathbf{x} according to the sign of $\beta_0 + \boldsymbol{\beta}'\mathbf{x}$, i.e., we classify $\mathbf{x} \sim \pi_1$ if $\beta_0 + \boldsymbol{\beta}'\mathbf{x} > 0$, and $\mathbf{x} \sim \pi_2$ if $\beta_0 + \boldsymbol{\beta}'\mathbf{x} < 0$. Then a point (y_i, \mathbf{x}_i) is misclassified iff $y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) < 0$.

The data sets are separable iff there exist β_0 and $\boldsymbol{\beta}$ such that $y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) > 0$ for all $i = 1, \dots, n$. The largest margin of separation can be obtained by solving the following problem⁵

$$\max_{\beta_0, \boldsymbol{\beta}, \|\boldsymbol{\beta}\|=1} c \quad (14.1)$$

$$\text{subject to } y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) \geq c, \quad i = 1, \dots, N. \quad (14.2)$$

The data is separable iff the optimal value of the above problem is positive. By making change of variables $c = 1/\|\boldsymbol{\beta}\|$, we can write the above problem as

$$\min_{\beta_0, \boldsymbol{\beta}} \|\boldsymbol{\beta}\|^2 \quad (14.3)$$

$$\text{subject to } y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) \geq 1, \quad i = 1, \dots, N. \quad (14.4)$$

Constraints (14.4) define a nonempty feasible set iff the data is separable. Problem (14.3) - (14.4) is a convex quadratic programming problem, and can be solved efficiently.

If the data sets (classes) overlap we can proceed in a similar way allowing some points to be on the wrong side of the margin. By introducing slack variables ξ_1, \dots, ξ_N we can modify the constraints $y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) \geq c$ as

$$y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) \geq c - \xi_i, \quad i = 1, \dots, N, \quad (14.5)$$

or

$$y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) \geq c(1 - \xi_i), \quad i = 1, \dots, N, \quad (14.6)$$

where $\xi_i \geq 0$, $i = 1, \dots, N$, and $\sum_{i=1}^N \xi_i \leq \text{const}$. Similar to (14.3)–(14.4), formulation (14.6) leads to the following optimization problem

$$\min_{\beta_0, \boldsymbol{\beta}, \boldsymbol{\xi}} \|\boldsymbol{\beta}\|^2 \quad (14.7)$$

$$\text{subject to } y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) \geq 1 - \xi_i, \quad i = 1, \dots, N, \quad (14.8)$$

$$\xi_i \geq 0, \quad i = 1, \dots, N, \quad (14.9)$$

$$\sum_{i=1}^N \xi_i \leq C, \quad (14.10)$$

where $C > 0$ is a chosen constant. The above problem (14.7) - (14.10) is a convex quadratic programming problem.

Recall that a point (y_i, \mathbf{x}_i) is misclassified iff $y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) < 0$. Therefore if $(\beta_0, \boldsymbol{\beta}, \boldsymbol{\xi})$ is a feasible point of the problem (14.7) - (14.10) and a point \mathbf{x}_i is misclassified, then $0 > y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) \geq 1 - \xi_i$, and hence $\xi_i > 1$. It follows that if C is smaller than the minimal number of possible misclassifications, then problem (14.7) - (14.10) does not have a feasible solution. On the other hand, for given β_0 and $\boldsymbol{\beta}$ consider the corresponding set of misclassified points. If C is equal to the number of misclassifications, then we can take $\xi_i = 1$ for every misclassified point and $\xi_i = 0$ for every classified point. This will give a feasible point of problem (14.7) - (14.10).

⁵The norm $\|\cdot\|$ here is the Euclidean norm.

We can look at the classification problem from the following point of view. Suppose that we want to find the hyperplane such that the number of misclassified points is minimal. That is, we would like to solve the following problem

$$\min_{\beta_0, \boldsymbol{\beta}} \sum_{i=1}^N \delta(-y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i)), \quad (14.11)$$

where $\delta(t) = 1$ if $t > 0$, and $\delta(t) = 0$ if $t \leq 0$. That is, for given β_0 and $\boldsymbol{\beta}$ the sum $\sum_{i=1}^N \delta(-y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i))$ is equal to the number of misclassified points.

Problem (14.11) is a difficult combinatorial problem. Note that $\delta(t) \leq [1 + t]_+$, where $[a]_+ = \max\{0, a\}$. Therefore we can approximate problem (14.11) by the following *convex* problem

$$\min_{\beta_0, \boldsymbol{\beta}} \sum_{i=1}^N [1 - y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i)]_+ + c\|\boldsymbol{\beta}\|^2. \quad (14.12)$$

Equivalently we can formulate problem (14.12) as

$$\min_{\beta_0, \boldsymbol{\beta}, \boldsymbol{\xi}} \quad \frac{1}{2}\|\boldsymbol{\beta}\|^2 + \gamma \sum_{i=1}^N \xi_i \quad (14.13)$$

$$\text{s.t.} \quad y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) \geq 1 - \xi_i, \quad i = 1, \dots, N, \quad (14.14)$$

$$\xi_i \geq 0, \quad i = 1, \dots, N, \quad (14.15)$$

where $\gamma = c^{-1}$.

The Lagrangian of the above problem is

$$L(\beta_0, \boldsymbol{\beta}, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \frac{1}{2}\|\boldsymbol{\beta}\|^2 + \gamma \sum_{i=1}^N \xi_i - \sum_{i=1}^N \lambda_i [y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) - (1 - \xi_i)] - \sum_{i=1}^N \mu_i \xi_i.$$

The Lagrangian dual of problem (14.13)–(14.15) is the problem

$$\max_{\boldsymbol{\lambda} \geq 0, \boldsymbol{\mu} \geq 0} \min_{\beta_0, \boldsymbol{\beta}, \boldsymbol{\xi} \geq 0} L(\beta_0, \boldsymbol{\beta}, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\mu}). \quad (14.16)$$

The corresponding Lagrangian-Wolfe dual is obtained by employing optimality conditions for the problem of minimization of $L(\beta_0, \boldsymbol{\beta}, \boldsymbol{\xi}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ in (14.16). That is, by setting derivatives of the Lagrangian to zero, with respect to $\boldsymbol{\beta}, \beta_0, \boldsymbol{\xi}$, we have

$$\boldsymbol{\beta} = \sum_{i=1}^N \lambda_i y_i \mathbf{x}_i \quad (14.17)$$

$$0 = \sum_{i=1}^N \lambda_i y_i \quad (14.18)$$

$$\lambda_i = \gamma - \mu_i, \quad i = 1, \dots, N, \quad (14.19)$$

By substituting these equations into the Lagrangian we obtain the Lagrangian-Wolfe dual:

$$\max_{\boldsymbol{\lambda}} \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{x}_i' \mathbf{x}_j \quad (14.20)$$

$$\text{s.t.} \quad 0 \leq \lambda_i \leq \gamma, \quad i = 1, \dots, N, \quad (14.21)$$

$$\sum_{i=1}^N \lambda_i y_i = 0. \quad (14.22)$$

We also have the following complementarity conditions for problem (14.13)– (14.15):

$$\lambda_i[y_i(\beta_0 + \boldsymbol{\beta}'\mathbf{x}_i) - (1 - \xi_i)] = 0, \quad i = 1, \dots, N, \quad (14.23)$$

$$\mu_i\xi_i = 0, \quad i = 1, \dots, N. \quad (14.24)$$

Given solution $\bar{\boldsymbol{\lambda}}$ of problem (14.20)–(14.22) the optimal $\boldsymbol{\beta}$ can be computed using equation (14.17), that is

$$\bar{\boldsymbol{\beta}} = \sum_{i=1}^N \bar{\lambda}_i y_i \mathbf{x}_i. \quad (14.25)$$

The complementarity conditions (14.24) mean that $\xi_i = 0$ if $\mu_i > 0$, and similarly for the complementarity conditions (14.23). By (14.19) we have that $\mu_i > 0$ if $\lambda_i < \gamma$. Therefore by using equation (14.23), for $0 < \bar{\lambda}_i < \gamma$ the optimal β_0 can be computed by solving $y_i f(\mathbf{x}_i) = 1$, where $f(\mathbf{x}) = \beta_0 + \boldsymbol{\beta}'\mathbf{x}$.

Suppose now that we want to make classification by using feature vectors $\mathbf{h}(\mathbf{x}_i)$, $i = 1, \dots, N$, where $\mathbf{h}(\cdot) = (h_1(\cdot), \dots, h_q(\cdot))' : \mathbb{R}^p \rightarrow \mathbb{R}^q$. We can approach this by solving the corresponding dual problem with replacing \mathbf{x}_i with $\mathbf{h}(\mathbf{x}_i)$, $i = 1, \dots, N$. That is the objective function in (14.20) is replaced by

$$\sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j \mathbf{h}(\mathbf{x}_i)' \mathbf{h}(\mathbf{x}_j). \quad (14.26)$$

Consequently, by using $\boldsymbol{\beta} = \sum_{i=1}^N \lambda_i y_i \mathbf{h}(\mathbf{x}_i)$ (see (14.25)), the classification is performed according to the sign of

$$f(\mathbf{x}) = \beta_0 + \boldsymbol{\beta}'\mathbf{h}(\mathbf{x}) = \beta_0 + \sum_{i=1}^N \lambda_i y_i \mathbf{h}(\mathbf{x}_i)' \mathbf{h}(\mathbf{x}). \quad (14.27)$$

Both expressions (14.26) and (14.27) are defined by the so-called kernel function

$$K(\mathbf{x}, \mathbf{z}) = \mathbf{h}(\mathbf{x})' \mathbf{h}(\mathbf{z}) = \sum_{s=1}^q h_s(\mathbf{x}) h_s(\mathbf{z}). \quad (14.28)$$

In terms of the kernel function the objective function (14.26) can be written as

$$\sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \lambda_i \lambda_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j), \quad (14.29)$$

and the classifier (14.27) as

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^N \lambda_i y_i K(\mathbf{x}, \mathbf{x}_i). \quad (14.30)$$

For example

$$K(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}'\mathbf{z})^2 = \left(1 + \sum_{i=1}^p x_i z_i\right)^2$$

defines a quadratic separation.

Kernel function should be symmetric, i.e., $K(\mathbf{x}, \mathbf{z}) = K(\mathbf{z}, \mathbf{x})$, and positive definite, i.e., for any $\mathbf{x}_1, \dots, \mathbf{x}_m$ the matrix $\mathbf{A} = [a_{ij}]$ with components $a_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$ should be positive semidefinite, or in other words $\sum_{i,j=1}^m \lambda_i \lambda_j K(\mathbf{x}_i, \mathbf{x}_j)$ should be nonnegative for any $\mathbf{x}_1, \dots, \mathbf{x}_m$ and $\lambda_1, \dots, \lambda_m$. Popular examples of kernels:

- Polynomial $K(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}'\mathbf{z})^d$.
- Radial basis $K(\mathbf{x}, \mathbf{z}) = \exp(-\gamma\|\mathbf{x} - \mathbf{z}\|^2)$, $\gamma > 0$.
- Hyperbolic tangent $K(\mathbf{x}, \mathbf{z}) = \tanh(c_1 + c_2\mathbf{x}'\mathbf{z})$, $c_1 < 0$, $c_2 > 0$, where $\tanh x = \frac{\sinh x}{\cosh x} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$, $\sinh x = -i \sin(ix) = \frac{e^x - e^{-x}}{2}$, $\cosh(x) = \cos(ix) = \frac{e^x + e^{-x}}{2}$.

15 Principal Components Analysis

Consider an $m \times 1$ random vector \mathbf{X} with $\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}]$ and $\boldsymbol{\Sigma} = \text{Cov}[\mathbf{X}]$. Let $\lambda_1 \geq \dots \geq \lambda_m$ be the eigenvalues and $\mathbf{e}_1, \dots, \mathbf{e}_m$ be corresponding eigenvectors of $\boldsymbol{\Sigma}$, i.e., $\boldsymbol{\Sigma}\mathbf{e}_i = \lambda_i\mathbf{e}_i$, $i = 1, \dots, m$. We assume⁶ that the eigenvectors are orthonormal, i.e., $\|\mathbf{e}_i\| = 1$, $i = 1, \dots, m$, and $\mathbf{e}_i'\mathbf{e}_j = 0$ for $i \neq j$. Recall that then (spectral decomposition)

$$\boldsymbol{\Sigma} = \mathbf{E}\boldsymbol{\Lambda}\mathbf{E}' = \lambda_1\mathbf{e}_1\mathbf{e}_1' + \dots + \lambda_m\mathbf{e}_m\mathbf{e}_m', \quad (15.1)$$

where $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_m)$ is diagonal matrix and $\mathbf{E} = [\mathbf{e}_1, \dots, \mathbf{e}_m]$ is orthogonal matrix.

Suppose that we want to find a linear combinations $\mathbf{w}'\mathbf{X} = w_1X_1 + \dots + w_mX_m$ with largest variance. That is we want to solve the problem

$$\max_{\|\mathbf{w}\|=1} \text{Var}(\mathbf{w}'\mathbf{X}). \quad (15.2)$$

Note that $\mathbf{w}'\mathbf{x} = \|\mathbf{w}\|\|\mathbf{x}\|\cos\theta$, where θ is the angle between vectors \mathbf{w} and \mathbf{x} . If $\|\mathbf{w}\| = 1$, then $\mathbf{w}'\mathbf{x} = \|\mathbf{x}\|\cos\theta$ is the orthogonal projection of vector \mathbf{x} onto the straight line in the direction of vector \mathbf{w} . Therefore problem (15.2) can be viewed as finding a direction such that projection of \mathbf{X} onto that direction has the largest variance.

We have that $\text{Var}(\mathbf{w}'\mathbf{X}) = \mathbf{w}'\boldsymbol{\Sigma}\mathbf{w}$ and by (15.1)

$$\mathbf{w}'\boldsymbol{\Sigma}\mathbf{w} = \mathbf{w}'\mathbf{E}\boldsymbol{\Lambda}\mathbf{E}'\mathbf{w} = \mathbf{v}'\boldsymbol{\Lambda}\mathbf{v} = \lambda_1v_1^2 + \dots + \lambda_mv_m^2, \quad (15.3)$$

where $\mathbf{v} = \mathbf{E}'\mathbf{w}$. Moreover, since matrix \mathbf{E} is orthogonal,

$$v_1^2 + \dots + v_m^2 = \mathbf{v}'\mathbf{v} = \mathbf{w}'\mathbf{E}\mathbf{E}'\mathbf{w} = \mathbf{w}'\mathbf{w} = 1.$$

That is, $\mathbf{v}'\boldsymbol{\Lambda}\mathbf{v}$ is a convex combination of eigenvalues λ_i . Thus $\mathbf{v}'\boldsymbol{\Lambda}\mathbf{v}$ is maximized when $\mathbf{v} = (1, 0, \dots, 0)'$. Since $\mathbf{w} = \mathbf{E}\mathbf{v}$, it follows that solution of problem (15.2) is given by the eigenvector \mathbf{e}_1 corresponding to the largest eigenvalue of matrix $\boldsymbol{\Sigma}$. Note that

$$\text{Var}(\mathbf{e}_1'\mathbf{X}) = \mathbf{e}_1'\boldsymbol{\Sigma}\mathbf{e}_1 = \lambda_1\mathbf{e}_1'\mathbf{e}_1 = \lambda_1.$$

Given the first principal component $Y_1 = \mathbf{e}_1'\mathbf{X}$, suppose that we want to find $Y_2 = \mathbf{w}'\mathbf{X}$, with $\|\mathbf{w}\| = 1$, such that $\text{Cov}(Y_1, Y_2) = 0$ and Y_2 has the largest possible variance. Since

$$\text{Cov}(Y_1, Y_2) = \mathbf{w}'\boldsymbol{\Sigma}\mathbf{e}_1 = \lambda_1\mathbf{w}'\mathbf{e}_1,$$

this means that we want to solve the problem

$$\max_{\|\mathbf{w}\|=1} \mathbf{w}'\boldsymbol{\Sigma}\mathbf{w} \quad \text{subject to } \mathbf{w}'\mathbf{e}_1 = 0. \quad (15.4)$$

⁶The norm in this section is the Euclidean norm.

Again we need to find \mathbf{v} which maximizes the right hand side of (15.3) and such that the sum of its squared components is one, and is orthogonal to vector $(1, 0, \dots, 0)'$, i.e., the first component of \mathbf{v} is zero. This is vector $(0, 1, 0, \dots, 0)'$, and hence solution of problem (15.4) is \mathbf{e}_2 .

And so on, variables $Y_i = \mathbf{e}_i' \mathbf{X}$, $i = 1, \dots, m$, are called principal components of the data vector \mathbf{X} . Note that $\text{Var}(Y_i) = \lambda_i$, $i = 1, \dots, m$, $\text{Cov}(Y_i, Y_j) = 0$ for $i \neq j$ and

$$\sum_{i=1}^m \text{Var}(Y_i) = \sum_{i=1}^m \text{Var}(X_i) = \sum_{i=1}^m \lambda_i = \text{tr}(\mathbf{\Sigma}).$$

Note also that vector \mathbf{Y} of principal components can be written as $\mathbf{Y} = \mathbf{E}' \mathbf{X}$. Multiplying both sides of this equation by \mathbf{E} , and since matrix \mathbf{E} is orthogonal, we obtain

$$\mathbf{X} = \mathbf{E}\mathbf{Y} = Y_1 \mathbf{e}_1 + \dots + Y_m \mathbf{e}_m. \quad (15.5)$$

That is, \mathbf{X} can be recovered from \mathbf{Y} if vectors $\mathbf{e}_1, \dots, \mathbf{e}_m$ are known. This can be used for approximation of \mathbf{X} by removing from the right hand side of (15.5) terms corresponding to small eigenvalues λ_i .

Note that principal components analysis is not scale invariant. That is, suppose we rescale components of \mathbf{X} say by changing units of measurements. So we change \mathbf{X} to $\mathbf{D}\mathbf{X}$, where \mathbf{D} is a diagonal matrix with positive diagonal elements representing change of scale. Then the covariance matrix $\mathbf{\Sigma}$ is changed to $\mathbf{D}\mathbf{\Sigma}\mathbf{D}$. The eigenvalues and eigenvectors of matrix $\mathbf{D}\mathbf{\Sigma}\mathbf{D}$ do not have a simple relation to the respective eigenvalues and eigenvectors of matrix $\mathbf{\Sigma}$.

The true (population) covariance matrix $\mathbf{\Sigma}$ is unknown. It is estimated by the sample covariance matrix

$$\mathbf{S} = (N - 1)^{-1} \sum_{i=1}^N (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})'. \quad (15.6)$$

Therefore the PCA usually performed on the sample covariance matrix \mathbf{S} , or because of the lack of scale invariance, on the sample correlation⁷ matrix. Let $\ell_1 \geq \dots \geq \ell_m$ be the eigenvalues and $\mathbf{q}_1, \dots, \mathbf{q}_m$ be corresponding orthonormal eigenvectors of \mathbf{S} , considered as estimates of the respective true eigenvalues and eigenvectors. What are statistical properties of these estimates? In order to apply Delta Theorem we need to compute derivatives of eigenvalues and eigenvectors considered as functions of symmetric matrices. We are going to discuss this next.

15.1 Derivatives of eigenvalues and eigenvectors

Consider the linear space of symmetric $m \times m$ matrices, denoted $\mathbb{S}^{m \times m}$. Consider $\mathbf{A} \in \mathbb{S}^{m \times m}$ and its eigenvalues $\lambda_1 \geq \dots \geq \lambda_m$ and the corresponding orthonormal eigenvectors $\mathbf{e}_1, \dots, \mathbf{e}_m$. Suppose that eigenvalue λ_i has *multiplicity one*, i.e., λ_i is different from the previous eigenvalue λ_{i-1} and the next eigenvalue λ_{i+1} . Then $\lambda_i(\cdot)$, considered as a function $\lambda_i : \mathbb{S}^{m \times m} \rightarrow \mathbb{R}$, is continuous at \mathbf{A} . Let us make small perturbations of elements of matrix \mathbf{A} by adding the differential $d\mathbf{A} \in \mathbb{S}^{m \times m}$. Then the eigenvalue equations for perturbed matrix are

$$(\mathbf{A} + d\mathbf{A})(\mathbf{e}_i + d\mathbf{e}_i) = (\lambda_i + d\lambda_i)(\mathbf{e}_i + d\mathbf{e}_i), \quad (15.7)$$

where $d\lambda_i$ and $d\mathbf{e}_i$ are the corresponding small changes in the eigenvalue and eigenvectors. Moreover we have that

$$(\mathbf{A} + d\mathbf{A})(\mathbf{e}_i + d\mathbf{e}_i) = \mathbf{A}\mathbf{e}_i + (d\mathbf{A})\mathbf{e}_i + \mathbf{A}(d\mathbf{e}_i) + (d\mathbf{A})d\mathbf{e}_i. \quad (15.8)$$

⁷The sample correlation matrix is obtained by scaling matrix \mathbf{DSD} such that its diagonal elements are equal one.

By disregarding the high order terms $(d\lambda_i)(de_i)$ and $(d\mathbf{A})(de_i)$ in (15.7) and (15.8), and since $\mathbf{A}e_i = \lambda_i e_i$, we can write

$$(d\mathbf{A})e_i + \mathbf{A}(de_i) = (d\lambda_i)e_i + \lambda_i(de_i). \quad (15.9)$$

Furthermore up to high order terms

$$(e_i + de_i)'(e_j + de_j) = (de_i)'e_j + e_i'(de_j) + e_i'e_j. \quad (15.10)$$

It follows that for $i = j$, since $(e_i + de_i)'(e_i + de_i) = e_i'e_i = 1$,

$$e_i'(de_i) = 0, \quad (15.11)$$

and for $i \neq j$, since $e_i'e_j = 0$ and $(e_i + de_i)'(e_j + de_j) = 0$,

$$(de_i)'e_j + e_i'(de_j) = 0. \quad (15.12)$$

Consequently by multiplying both sides of (15.9) by e_i' and noting that $e_i'e_i = 1$, $e_i'(de_i) = 0$ and $e_i'\mathbf{A}(de_i) = \lambda_i e_i'(de_i) = 0$, we obtain

$$d\lambda_i = e_i'(d\mathbf{A})e_i. \quad (15.13)$$

It is also possible to write (15.13) as

$$d\lambda_i = \text{tr}(e_i e_i'(d\mathbf{A})). \quad (15.14)$$

Equation (15.13) (equation (15.14)) gives an expression for the linear approximation of the eigenvalue λ_i for small perturbations $d\mathbf{A}$ of matrix \mathbf{A} , i.e.,

$$\lambda_i(\mathbf{A} + d\mathbf{A}) = \lambda_i(\mathbf{A}) + e_i'(d\mathbf{A})e_i + o(\|d\mathbf{A}\|). \quad (15.15)$$

The assumption that the eigenvalue is simple is essential in the above derivations.

Now let us compute de_i . Note that since it is assumed that the eigenvalue λ_i is simple and $\|e_i\| = 1$, the eigenvector e_i of \mathbf{A} is defined uniquely up to sign change from e_i to $-e_i$. Since eigenvectors e_1, \dots, e_m are orthonormal, they form a basis and hence we can write de_i as linear combination $de_i = c_1 e_1 + \dots + c_m e_m$ with $c_j = e_j' de_i$, $j = 1, \dots, m$. For $i \neq j$ we have by (15.9) and since $e_j' e_i = 0$ that

$$e_j'(d\mathbf{A})e_i + e_j'\mathbf{A}(de_i) = \lambda_i e_j'(de_i), \quad (15.16)$$

and since $e_j'\mathbf{A}(de_i) = \lambda_j e_j'(de_i)$ it follows that

$$e_j'(d\mathbf{A})e_i = (\lambda_i - \lambda_j) e_j'(de_i). \quad (15.17)$$

This implies that

$$c_j = (\lambda_i - \lambda_j)^{-1} e_j'(d\mathbf{A})e_i, \quad j \neq i. \quad (15.18)$$

For $j = i$ we have $c_i = e_i'(de_i) = 0$. We obtain the following formula for the differential of e_i :

$$de_i = \sum_{\substack{j=1 \\ j \neq i}}^m \left[\frac{e_j'(d\mathbf{A})e_i}{\lambda_i - \lambda_j} \right] e_j. \quad (15.19)$$

That is, for small perturbations $d\mathbf{A}$ of matrix \mathbf{A} ,

$$e_i(\mathbf{A} + d\mathbf{A}) = \sum_{\substack{j=1 \\ j \neq i}}^m \left[\frac{e_j'(d\mathbf{A})e_i}{\lambda_i - \lambda_j} \right] e_j + o(\|d\mathbf{A}\|). \quad (15.20)$$

15.2 Elements of matrix calculus

Kronecker product of matrices $\mathbf{A} = [a_{ij}]$ and $\mathbf{B} = [b_{ij}]$, of respective orders $p \times q$ and $r \times s$, is the $pr \times qs$ matrix

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1q}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2q}\mathbf{B} \\ \cdots & \cdots & \cdots & \cdots \\ a_{p1}\mathbf{B} & a_{p2}\mathbf{B} & \cdots & a_{pq}\mathbf{B} \end{bmatrix}.$$

Vec-operator of $p \times q$ matrix \mathbf{A} is $pq \times 1$ vector

$$\text{vec}(\mathbf{A}) = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_q \end{bmatrix},$$

where $\mathbf{a}_1, \dots, \mathbf{a}_q$ are columns of \mathbf{A} .

Note the following matrix identities

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD}) \quad (15.21)$$

and

$$\text{vec}(\mathbf{BXC}) = (\mathbf{C}' \otimes \mathbf{B})\text{vec}(\mathbf{X}), \quad (15.22)$$

where $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}, \mathbf{X}$ are matrices of appropriate order. Also for matrices \mathbf{A} and \mathbf{B} of the same order $p \times q$, and vectors $\mathbf{a} = \text{vec}(\mathbf{A})$ and $\mathbf{b} = \text{vec}(\mathbf{B})$,

$$\text{tr}(\mathbf{A}'\mathbf{B}) = \sum_{i,j} a_{ij}b_{ij} = \mathbf{a}'\mathbf{b}. \quad (15.23)$$

Let $\mathbf{X}_1, \dots, \mathbf{X}_N$ be an iid sample of realizations of random vector $\mathbf{X} = (X_1, \dots, X_m)'$. Assume that the distribution of \mathbf{X} has finite fourth order moments. Let $\mathbf{s} = \text{vec}(\mathbf{S})$ and $\boldsymbol{\sigma}_0 = \text{vec}(\boldsymbol{\Sigma}_0)$, where $\boldsymbol{\Sigma}_0 = [\sigma_{ij}]$ is the population covariance matrix. Then by the CLT, $N^{1/2}(\mathbf{s} - \boldsymbol{\sigma}_0)$ converges in distribution to normal with zero mean vector and a covariance matrix $\boldsymbol{\Gamma}$ of order $m^2 \times m^2$. Note that since matrices \mathbf{S} and $\boldsymbol{\Sigma}_0$ are symmetric, vectors \mathbf{s} and $\boldsymbol{\sigma}_0$ have not more than $m(m+1)/2$ different elements, therefore $\text{rank}(\boldsymbol{\Gamma}) \leq m(m+1)/2$. The typical element of matrix $\boldsymbol{\Gamma}$ is

$$\begin{aligned} [\boldsymbol{\Gamma}]_{ij,kl} &= \mathbb{E}\{[(X_i - \mu_i)(X_j - \mu_j) - \sigma_{ij}][(X_k - \mu_k)(X_\ell - \mu_\ell) - \sigma_{kl}]\} \\ &= \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)(X_k - \mu_k)(X_\ell - \mu_\ell)] - \sigma_{ij}\sigma_{kl}. \end{aligned}$$

In particular if \mathbf{X} has normal distribution, then

$$\mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)(X_k - \mu_k)(X_\ell - \mu_\ell)] = \sigma_{ij}\sigma_{kl} + \sigma_{ik}\sigma_{jl} + \sigma_{il}\sigma_{jk},$$

and hence

$$[\boldsymbol{\Gamma}]_{ij,kl} = \sigma_{ik}\sigma_{jl} + \sigma_{il}\sigma_{jk}. \quad (15.24)$$

In a matrix form equations (15.24) can be written as

$$\boldsymbol{\Gamma} = 2\mathbf{M}_m(\boldsymbol{\Sigma}_0 \otimes \boldsymbol{\Sigma}_0), \quad (15.25)$$

where \mathbf{M}_m is the $m^2 \times m^2$ matrix given by

$$\mathbf{M}_m = \frac{1}{2} \left[\mathbf{I}_{m^2} + \sum_{i,j=1}^m (\mathbf{H}_{ij} \otimes \mathbf{H}'_{ij}) \right],$$

with \mathbf{H}_{ij} being $m \times m$ matrix with $h_{ij} = 1$ and all other elements zero. The matrix \mathbf{M}_m has the following properties: (i) $\text{rank}(\mathbf{M}_m) = m(m+1)/2$, (ii) $\mathbf{M}_m^2 = \mathbf{M}_m$, (iii) for any symmetric matrix $\mathbf{\Sigma}$,

$$\mathbf{M}_m(\mathbf{\Sigma} \otimes \mathbf{\Sigma}) = (\mathbf{\Sigma} \otimes \mathbf{\Sigma})\mathbf{M}_m \text{ and } \mathbf{M}_m \text{vec}(\mathbf{\Sigma}) = \text{vec}(\mathbf{\Sigma}).$$

It follows that

$$\mathbf{\Gamma} = 2\mathbf{M}_m(\mathbf{\Sigma}_0 \otimes \mathbf{\Sigma}_0)\mathbf{M}_m. \quad (15.26)$$

15.3 Asymptotics of PCA

Let $\mathbf{X}_1, \dots, \mathbf{X}_N$ be an iid sample from $\mathcal{N}_m(\boldsymbol{\mu}, \mathbf{\Sigma})$ and \mathbf{S} be the corresponding sample covariance matrix. Let $\lambda_1 \geq \dots \geq \lambda_m$ be the eigenvalues and $\mathbf{e}_1, \dots, \mathbf{e}_m$ be a corresponding set of orthonormal eigenvectors of $\mathbf{\Sigma}$, and $\ell_1 \geq \dots \geq \ell_m$ be the eigenvalues and $\mathbf{q}_1, \dots, \mathbf{q}_m$ be a corresponding set of orthonormal eigenvectors of \mathbf{S} .

Suppose that λ_i has multiplicity one. Let us show that $N^{1/2}(\ell_i - \lambda_i)$ and $N^{1/2}(\mathbf{q}_i - \mathbf{e}_i)$ are asymptotically normal (with mean zero) and asymptotically independent of each other, and that the asymptotic variance of $N^{1/2}(\ell_i - \lambda_i)$ is $2\lambda_i^2$ and the asymptotic covariance matrix of $N^{1/2}(\mathbf{q}_i - \mathbf{e}_i)$ is

$$\sum_{j=1, j \neq i}^m \frac{\lambda_i \lambda_j}{(\lambda_i - \lambda_j)^2} \mathbf{e}_j \mathbf{e}'_j. \quad (15.27)$$

By the Delta Theorem and (15.15) we have that

$$N^{1/2}(\ell_i - \lambda_i) = \mathbf{e}'_i [N^{1/2}(\mathbf{S} - \mathbf{\Sigma})] \mathbf{e}_i + o_p(1),$$

and hence $N^{1/2}(\ell_i - \lambda_i)$ converges in distribution to $\mathcal{N}(0, \sigma^2)$, where σ^2 can be calculated as follows. We have

$$\mathbf{e}'_i [N^{1/2}(\mathbf{S} - \mathbf{\Sigma})] \mathbf{e}_i = \text{tr} \left[N^{1/2}(\mathbf{S} - \mathbf{\Sigma}) \mathbf{e}_i \mathbf{e}'_i \right] = [\text{vec}(\mathbf{e}_i \mathbf{e}'_i)]' [N^{1/2}(\mathbf{s} - \boldsymbol{\sigma})]$$

and hence

$$\sigma^2 = [\text{vec}(\mathbf{e}_i \mathbf{e}'_i)]' \mathbf{\Gamma} [\text{vec}(\mathbf{e}_i \mathbf{e}'_i)] = 2[\text{vec}(\mathbf{e}_i \mathbf{e}'_i)]' \mathbf{M}_m(\mathbf{\Sigma} \otimes \mathbf{\Sigma}) \mathbf{M}_m [\text{vec}(\mathbf{e}_i \mathbf{e}'_i)].$$

Moreover, $\mathbf{M}_m [\text{vec}(\mathbf{e}_i \mathbf{e}'_i)] = \text{vec}(\mathbf{e}_i \mathbf{e}'_i)$ and

$$[\text{vec}(\mathbf{e}_i \mathbf{e}'_i)]' (\mathbf{\Sigma} \otimes \mathbf{\Sigma}) [\text{vec}(\mathbf{e}_i \mathbf{e}'_i)] = \text{tr}[(\mathbf{e}_i \mathbf{e}'_i) \mathbf{\Sigma} (\mathbf{e}_i \mathbf{e}'_i) \mathbf{\Sigma}] = (\mathbf{e}'_i \mathbf{\Sigma} \mathbf{e}_i) (\mathbf{e}'_i \mathbf{\Sigma} \mathbf{e}_i) = \lambda_i^2.$$

Similarly, by (15.20)

$$N^{1/2}(\mathbf{q}_i - \mathbf{e}_i) = \sum_{j \neq i} a_j \mathbf{e}_j + o_p(1),$$

where

$$a_j = \frac{\mathbf{e}'_j [N^{1/2}(\mathbf{S} - \mathbf{\Sigma})] \mathbf{e}_i}{\lambda_i - \lambda_j} = (\lambda_i - \lambda_j)^{-1} [\text{vec}(\mathbf{e}_i \mathbf{e}'_i)]' [N^{1/2}(\mathbf{s} - \boldsymbol{\sigma})].$$

The asymptotic covariance between a_j and a_k (for $j \neq i$ and $k \neq i$) is

$$\frac{[\text{vec}(\mathbf{e}_i \mathbf{e}'_j)]' \mathbf{\Gamma} [\text{vec}(\mathbf{e}_i \mathbf{e}'_k)]}{(\lambda_i - \lambda_j)^2} = \frac{2 \text{tr}[(\mathbf{e}_i \mathbf{e}'_j) \mathbf{M}_m \mathbf{\Sigma} (\mathbf{e}_i \mathbf{e}'_k) \mathbf{M}_m \mathbf{\Sigma}]}{(\lambda_i - \lambda_j)^2}. \quad (15.28)$$

Also $\mathbf{M}_m(\mathbf{e}_i \mathbf{e}'_j) = \frac{1}{2}[(\mathbf{e}_i \mathbf{e}'_j) + (\mathbf{e}_j \mathbf{e}'_i)]$ and $\mathbf{M}_m(\mathbf{e}_i \mathbf{e}'_k) = \frac{1}{2}[(\mathbf{e}_i \mathbf{e}'_k) + (\mathbf{e}_k \mathbf{e}'_i)]$. It follows that the right hand side of (15.28) is equal to

$$\frac{\text{tr}[(\mathbf{e}_i \mathbf{e}'_j) + (\mathbf{e}_j \mathbf{e}'_i)] \mathbf{\Sigma} [(\mathbf{e}_i \mathbf{e}'_k) + (\mathbf{e}_k \mathbf{e}'_i)] \mathbf{\Sigma}}{2(\lambda_i - \lambda_j)^2}. \quad (15.29)$$

Moreover, we have that $\mathbf{e}'_j \mathbf{\Sigma} \mathbf{e}_k = \lambda_j \mathbf{e}'_j \mathbf{e}_k$ equals 0 if $j \neq k$, and λ_j if $j = k$. Therefore the expression in (15.29) equals 0 if $j \neq k$, and $\lambda_i \lambda_j / (\lambda_i - \lambda_j)^2$ if $j = k$. We obtain that the asymptotic covariance matrix of $a_j \mathbf{e}_j$ is $\frac{\lambda_i \lambda_j}{(\lambda_i - \lambda_j)^2} \mathbf{e}_j \mathbf{e}'_j$, and $a_j \mathbf{e}_j$ is asymptotically uncorrelated with $a_k \mathbf{e}_k$ for $j \neq k$. Formula (15.27) follows.

Finally, the asymptotic covariance between $n^{1/2}(\ell_i - \lambda_i)$ and a_j , $j \neq i$, is proportional to

$$[\text{vec}(\mathbf{e}_i \mathbf{e}'_i)] \mathbf{\Gamma} [\text{vec}(\mathbf{e}_i \mathbf{e}'_j)] = \text{tr}[(\mathbf{e}_i \mathbf{e}'_i) \mathbf{\Sigma} (\mathbf{e}_i \mathbf{e}'_j + \mathbf{e}_j \mathbf{e}'_i) \mathbf{\Sigma}] = 0,$$

and hence $N^{1/2}(\ell_i - \lambda_i)$ and $N^{1/2}(\mathbf{q}_i - \mathbf{e}_i)$ are asymptotically independent.

15.4 Singular value decomposition

Let \mathbf{X} be an $m \times n$ matrix of rank r (note that $r \leq \min\{m, n\}$). Its singular value decomposition is

$$\mathbf{X} = \mathbf{V} \mathbf{D} \mathbf{W}' = \sigma_1 \mathbf{v}_1 \mathbf{w}'_1 + \dots + \sigma_r \mathbf{v}_r \mathbf{w}'_r,$$

where $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_r]$ and $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_r]$ are matrices of order $m \times r$ and $n \times r$, respectively, such that $\mathbf{V}' \mathbf{V} = \mathbf{I}_r$ and $\mathbf{W}' \mathbf{W} = \mathbf{I}_r$, and $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_r)$ with $\sigma_1 \geq \dots \geq \sigma_r > 0$. Note that

$$\mathbf{X} \mathbf{X}' = \mathbf{V} \mathbf{D} \mathbf{W}' \mathbf{W} \mathbf{D} \mathbf{V} = \mathbf{V} \mathbf{D}^2 \mathbf{V}',$$

i.e., $\mathbf{V} \mathbf{D}^2 \mathbf{V}'$ is the spectral decomposition of the (symmetric positive semidefinite) $m \times m$ matrix $\mathbf{X} \mathbf{X}'$. It follows that σ_i^2 , $i = 1, \dots, r$, are the nonzero eigenvalues of $\mathbf{X} \mathbf{X}'$. Similarly $\mathbf{W} \mathbf{D}^2 \mathbf{W}'$ is the spectral decomposition of the $n \times n$ matrix $\mathbf{X}' \mathbf{X}$ with the same nonzero eigenvalues σ_i^2 , $i = 1, \dots, r$.

For $s < r$ consider the (truncated) matrix $\mathbf{X}_s = \mathbf{V}_s \mathbf{D}_s \mathbf{W}'_s$, where $\mathbf{V}_s = [\mathbf{v}_1, \dots, \mathbf{v}_s]$, $\mathbf{W}_s = [\mathbf{w}_1, \dots, \mathbf{w}_s]$ and $\mathbf{D}_s = \text{diag}(\sigma_1, \dots, \sigma_s)$, i.e.,

$$\mathbf{X}_s = \sigma_1 \mathbf{v}_1 \mathbf{w}'_1 + \dots + \sigma_s \mathbf{v}_s \mathbf{w}'_s.$$

The matrix \mathbf{X}_s is the nearest matrix of rank s to the original matrix \mathbf{X} , in the sense of the difference between the two having the smallest possible Frobenius norm (Eckart-Young theorem). That is, solution of the minimization problem

$$\min_{\mathbf{Z} \in \mathbb{R}^{m \times n}} \|\mathbf{X} - \mathbf{Z}\|_F \quad \text{subject to } \text{rank}(\mathbf{Z}) \leq s$$

is $\bar{\mathbf{Z}} = \mathbf{X}_s$. Frobenius norm of a matrix \mathbf{A} is

$$\|\mathbf{A}\|_F = \sqrt{\text{tr}(\mathbf{A} \mathbf{A}')} = \sqrt{\text{tr}(\mathbf{A}' \mathbf{A})} = \sqrt{\sum_{i,j} a_{ij}^2}.$$

16 Factor analysis model

Consider an $m \times 1$ random vector \mathbf{X} (of measurements) with $\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}]$ and $\text{Cov}(\mathbf{X}) = \boldsymbol{\Sigma}$. The factor analysis model assumes that

$$\mathbf{X} = \boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{f} + \boldsymbol{\varepsilon}, \quad (16.1)$$

where $\boldsymbol{\Lambda}$ is an $m \times k$ matrix (of factor loadings), \mathbf{f} is a $k \times 1$ random vector (of factors) and $\boldsymbol{\varepsilon}$ is an $m \times 1$ random vector (errors). It is assumed that: (i) $\mathbb{E}[\mathbf{f}] = \mathbf{0}$ and $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$, (ii) the errors are uncorrelated, i.e., $\text{Cov}(\boldsymbol{\varepsilon})$ is diagonal, (iii) the factors and errors are uncorrelated, i.e., $\mathbb{E}[\mathbf{f}\boldsymbol{\varepsilon}'] = \mathbf{0}$.

It follows then that

$$\boldsymbol{\Sigma} = \mathbb{E}[(\boldsymbol{\Lambda}\mathbf{f} + \boldsymbol{\varepsilon})(\boldsymbol{\Lambda}\mathbf{f} + \boldsymbol{\varepsilon})'] = \boldsymbol{\Lambda}\boldsymbol{\Phi}\boldsymbol{\Lambda}' + \boldsymbol{\Psi}, \quad (16.2)$$

where $\boldsymbol{\Phi} = \text{Cov}(\mathbf{f})$ and $\boldsymbol{\Psi} = \text{Cov}(\boldsymbol{\varepsilon})$. Since it is assumed that the errors are uncorrelated, matrix $\boldsymbol{\Psi}$ is diagonal. Matrices $\boldsymbol{\Phi}$ and $\boldsymbol{\Psi}$ are covariance matrices and hence are positive semidefinite. Since matrix $\boldsymbol{\Psi}$ is diagonal, it is positive semidefinite iff all its diagonal elements are nonnegative. Often it is assumed that $\boldsymbol{\Phi} = \mathbf{I}_k$, i.e., the factors are standardized. Then the model becomes

$$\boldsymbol{\Sigma} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}' + \boldsymbol{\Psi}. \quad (16.3)$$

Note that $\text{rank}(\boldsymbol{\Sigma} - \boldsymbol{\Psi}) = \text{rank}(\boldsymbol{\Lambda}) \leq k$. Note also that if \mathbf{T} is a $k \times k$ orthogonal matrix, then $(\boldsymbol{\Lambda}\mathbf{T})(\boldsymbol{\Lambda}\mathbf{T})' = \boldsymbol{\Lambda}\mathbf{T}\mathbf{T}'\boldsymbol{\Lambda}' = \boldsymbol{\Lambda}\boldsymbol{\Lambda}'$. Therefore the right hand side of (16.3) is defined up to change of $\boldsymbol{\Lambda}$ to $\boldsymbol{\Lambda}\mathbf{T}$. This can be viewed as rotation of the row vectors of matrix $\boldsymbol{\Lambda}$ by orthogonal matrix \mathbf{T} .

There is a certain similarity between Factor Analysis (FA) and PCA. Both try to explain covariances between components of the response vector \mathbf{X} by a smaller number of factors. But there are also essential differences, FA is a model and, unlike PCA, is scale invariant. That is, if \mathbf{D} is a diagonal matrix with positive diagonal elements, then rescaling \mathbf{X} to $\mathbf{D}\mathbf{X}$ results in rescaling $\boldsymbol{\Sigma}$ to $\mathbf{D}\boldsymbol{\Sigma}\mathbf{D}$, $\boldsymbol{\Lambda}$ to $\mathbf{D}\boldsymbol{\Lambda}$ and $\boldsymbol{\Psi}$ to $\mathbf{D}^2\boldsymbol{\Psi}$. It is possible to develop a statistical inference of the FA model (below).

Given data (sample) $\mathbf{X}_1, \dots, \mathbf{X}_N$ of observations (realizations) of \mathbf{X} , FA is performed on the sample covariance matrix \mathbf{S} . That is, \mathbf{S} is approximated by matrix of the form $\hat{\boldsymbol{\Sigma}} = \hat{\boldsymbol{\Lambda}}\hat{\boldsymbol{\Lambda}}' + \hat{\boldsymbol{\Psi}}$, where $\hat{\boldsymbol{\Lambda}}$ is an $m \times k$ matrix and $\hat{\boldsymbol{\Psi}}$ is a diagonal matrix with nonnegative diagonal elements. If it is assumed that the sample has normal distribution, i.e., $\mathbf{X}_i \sim \mathcal{N}_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, it is possible to proceed to the corresponding statistical inference. Let us show that the ML estimators of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are $\hat{\boldsymbol{\mu}} = \bar{\mathbf{X}}$ and $\hat{\boldsymbol{\Sigma}} = N^{-1} \sum_{i=1}^N (\mathbf{X}_i - \bar{\mathbf{X}})'(\mathbf{X}_i - \bar{\mathbf{X}})$. Note that $\hat{\boldsymbol{\Sigma}} = \frac{N-1}{N}\mathbf{S}$.

The likelihood function is

$$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{i=1}^N \frac{1}{(2\pi)^{m/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -(\mathbf{X}_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{X}_i - \boldsymbol{\mu}) / 2 \right\}.$$

Up to a constant independent of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ we can write logarithm of the likelihood function as

$$\begin{aligned}
\log L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= -\frac{1}{2}N \log |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{i=1}^N (\mathbf{X}_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{X}_i - \boldsymbol{\mu}) \\
&= -\frac{1}{2}N \log |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{i=1}^N \text{tr} [\boldsymbol{\Sigma}^{-1} (\mathbf{X}_i - \boldsymbol{\mu})(\mathbf{X}_i - \boldsymbol{\mu})'] \\
&= -\frac{1}{2}N \log |\boldsymbol{\Sigma}| - \frac{1}{2} \text{tr} \boldsymbol{\Sigma}^{-1} \left[\sum_{i=1}^N (\mathbf{X}_i - \boldsymbol{\mu})(\mathbf{X}_i - \boldsymbol{\mu})' \right] \\
&= -\frac{1}{2}N \log |\boldsymbol{\Sigma}| - \frac{1}{2} \text{tr} \left\{ \boldsymbol{\Sigma}^{-1} \left[\sum_{i=1}^N (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})' + N(\bar{\mathbf{X}} - \boldsymbol{\mu})(\bar{\mathbf{X}} - \boldsymbol{\mu})' \right] \right\} \\
&= -\frac{1}{2}N \log |\boldsymbol{\Sigma}| - \frac{1}{2} \text{tr} [\boldsymbol{\Sigma}^{-1} \mathbf{A}] - \frac{1}{2}N(\bar{\mathbf{X}} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}).
\end{aligned}$$

where

$$\mathbf{A} := \sum_{i=1}^N (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})' = (N-1)\mathbf{S}.$$

That is

$$\log L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{1}{2}N \log |\boldsymbol{\Sigma}| - \frac{1}{2}(N-1) \text{tr} [\boldsymbol{\Sigma}^{-1} \mathbf{S}] - \frac{1}{2}N(\bar{\mathbf{X}} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}). \quad (16.4)$$

Since $\boldsymbol{\Sigma}^{-1}$ is positive definite, we have that $(\bar{\mathbf{X}} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}) \geq 0$ and its minimum of zero is attained for $\boldsymbol{\mu} = \bar{\mathbf{X}}$. It follows that $\bar{\mathbf{X}}$ is the ML estimator of $\boldsymbol{\mu}$. Now in order to find the ML estimator of $\boldsymbol{\Sigma}$ we need to minimize $N \log |\boldsymbol{\Sigma}| + \text{tr} [\boldsymbol{\Sigma}^{-1} \mathbf{A}]$ over positive definite matrices $\boldsymbol{\Sigma}$. Let $\lambda_1, \dots, \lambda_m$ be eigenvalues of $\boldsymbol{\Sigma}^{-1} \mathbf{A}$ (note that since matrices $\boldsymbol{\Sigma}$ and \mathbf{A} are positive definite, matrix $\boldsymbol{\Sigma}^{-1} \mathbf{A}$ has positive real valued eigenvalues, see section 1). Then

$$N \log |\boldsymbol{\Sigma}| + \text{tr} [\boldsymbol{\Sigma}^{-1} \mathbf{A}] = N \log |\boldsymbol{\Sigma} \mathbf{A}^{-1}| + \text{tr} [\boldsymbol{\Sigma}^{-1} \mathbf{A}] + N \log |\mathbf{A}| = \sum_{i=1}^m (\lambda_i - N \log \lambda_i) + N \log |\mathbf{A}|.$$

Note that function $f(\lambda) = \lambda - N \log \lambda$ is convex and has unique minimizer $\lambda = N$. It follows that the minimum is attained when all eigenvalues $\lambda_i = N$, that is $\hat{\boldsymbol{\Sigma}}^{-1} \mathbf{A} = N \mathbf{I}_m$. It follows that $\hat{\boldsymbol{\Sigma}} = N^{-1} \mathbf{A} = \frac{N-1}{N} \mathbf{S}$. \square

Assuming that the sample is from normally distributed population, by (16.4) and since the MLE of $\boldsymbol{\mu}$ is $\bar{\mathbf{X}}$, the MLE of parameters $\boldsymbol{\Lambda}$ and $\boldsymbol{\Psi}$ of the FA model are obtained by solving the problem

$$\min_{\boldsymbol{\Lambda}, \boldsymbol{\Psi} \geq 0} \log |\boldsymbol{\Lambda} \boldsymbol{\Lambda}' + \boldsymbol{\Psi}| + \frac{N-1}{N} \text{tr} [(\boldsymbol{\Lambda} \boldsymbol{\Lambda}' + \boldsymbol{\Psi})^{-1} \mathbf{S}],$$

where matrix $\boldsymbol{\Psi}$ is diagonal (by writing $\boldsymbol{\Psi} \geq 0$ we mean that diagonal elements of $\boldsymbol{\Psi}$ are nonnegative).

An important question in FA is how many factors should be in the model. The LRT statistic for testing FA model (16.3), with k factors, is

$$2 \log \lambda = N \min_{\boldsymbol{\Lambda}, \boldsymbol{\Psi} \geq 0} \left\{ \log |\boldsymbol{\Lambda} \boldsymbol{\Lambda}' + \boldsymbol{\Psi}| - \log |\hat{\boldsymbol{\Sigma}}| + \text{tr} [(\boldsymbol{\Lambda} \boldsymbol{\Lambda}' + \boldsymbol{\Psi})^{-1} \hat{\boldsymbol{\Sigma}}] - m \right\}$$

where $\hat{\boldsymbol{\Sigma}} = \frac{N-1}{N} \mathbf{S}$ is the unrestricted MLE of $\boldsymbol{\Sigma}$. Under H_0 of the FA model with k factors, the statistic $2 \log \lambda$ asymptotically has χ_ν^2 distribution with $\nu = m(m+1)/2 - m(k+1) + k(k-1)/2$

degrees of freedom. In calculation of the degrees of freedom, $m(m+1)/2$ is the number of nonduplicated elements of the covariance matrix, $mk+m$ is the number of estimated parameters and the last term $k(k-1)/2$ is the correction because of the possible rotation of the factor loadings matrix by $k \times k$ orthogonal matrix. Consequently H_0 hypothesis of k factors is rejected if the statistic $2 \log \lambda$ is larger than critical value of the χ^2_ν distribution.

The above statistical inference is based on the assumption that the population has a normal distribution. In various applications this assumption can be questionable. Also if the sample size n is large, this procedure tends to reject H_0 even if the FA model gives a reasonable approximation of the sample covariance matrix. Various indexes of fit, with questionable justifications, were suggested in the literature trying to resolve the question of ‘correct’ number of factors.

17 Kernel PCA

Given data (sample) $\mathbf{X}_1, \dots, \mathbf{X}_N$, suppose that we want to represent data in terms of vectors $\mathbf{Z}_i = \mathbf{h}(\mathbf{X}_i)$, $i = 1, \dots, N$, where $\mathbf{h}(\cdot) = (h_1(\cdot), \dots, h_q(\cdot))' : \mathbb{R}^m \rightarrow \mathbb{R}^q$ is a (nonlinear) mapping. Suppose for the moment that $\bar{\mathbf{Z}} = N^{-1} \sum_{i=1}^N \mathbf{Z}_i = N^{-1} \sum_{i=1}^N \mathbf{h}(\mathbf{X}_i)$ is $\mathbf{0}$. Consider the corresponding estimator of the covariance matrix in the new feature space

$$\mathbf{C} = N^{-1} \sum_{i=1}^N \mathbf{Z}_i \mathbf{Z}_i' = N^{-1} \sum_{i=1}^N \mathbf{h}(\mathbf{X}_i) \mathbf{h}(\mathbf{X}_i)'$$

Let $\lambda_1 \geq \dots \geq \lambda_q$ be eigenvalues and $\mathbf{e}_1, \dots, \mathbf{e}_q$ be corresponding orthonormal eigenvectors of the $q \times q$ matrix \mathbf{C} , i.e., $\mathbf{C} \mathbf{e}_s = \lambda_s \mathbf{e}_s$, $s = 1, \dots, q$. We have that

$$\lambda_s \mathbf{e}_s = \mathbf{C} \mathbf{e}_s = N^{-1} \sum_{i=1}^N \mathbf{Z}_i \mathbf{Z}_i' \mathbf{e}_s,$$

and hence (for $\lambda_s \neq 0$)

$$\mathbf{e}_s = \frac{1}{\lambda_s N} \sum_{i=1}^N \alpha_{is} \mathbf{Z}_i, \quad (17.1)$$

where $\alpha_{is} = \mathbf{Z}_i' \mathbf{e}_s$, $s = 1, \dots, q$, $i = 1, \dots, N$. It follows by (17.1) that

$$\alpha_{is} = \frac{1}{\lambda_s N} \mathbf{Z}_i' \left(\sum_{j=1}^N \alpha_{js} \mathbf{Z}_j \right) = \frac{1}{\lambda_s N} \sum_{j=1}^N \alpha_{js} \mathbf{Z}_i' \mathbf{Z}_j = \frac{1}{\lambda_s N} \sum_{j=1}^N \alpha_{js} \mathbf{h}(\mathbf{X}_i)' \mathbf{h}(\mathbf{X}_j). \quad (17.2)$$

Consider kernel function (compare with (14.28)) $K(\mathbf{x}, \mathbf{z}) = \mathbf{h}(\mathbf{x})' \mathbf{h}(\mathbf{z})$. In terms of the kernel function equation (17.2) can be written as

$$\sum_{j=1}^N \alpha_{js} K(\mathbf{X}_i, \mathbf{X}_j) = \lambda_s N \alpha_{is}. \quad (17.3)$$

Consider $N \times N$ matrix \mathbf{K} with components $\mathbf{K}_{ij} = K(\mathbf{X}_i, \mathbf{X}_j)$, $i, j = 1, \dots, N$. Equation (17.3) can be written as

$$\mathbf{K} \boldsymbol{\alpha}_s = \lambda_s N \boldsymbol{\alpha}_s, \quad s = 1, \dots, q, \quad (17.4)$$

where $\boldsymbol{\alpha}_s = (\alpha_{1s}, \dots, \alpha_{Ns})'$. That is, $\boldsymbol{\alpha}_s$ are eigenvectors of matrix \mathbf{K} . These eigenvectors can be normalized as follows

$$1 = \mathbf{e}_s' \mathbf{e}_s = \frac{1}{\lambda_s^2 N^2} \left(\sum_{i=1}^N \alpha_{is} \mathbf{Z}_i' \right) \left(\sum_{j=1}^N \alpha_{js} \mathbf{Z}_j \right) = \frac{1}{\lambda_s^2 N^2} \sum_{i,j=1}^N \alpha_{is} \alpha_{js} \mathbf{Z}_i' \mathbf{Z}_j = \frac{1}{\lambda_s^2 N^2} \sum_{i,j=1}^N \alpha_{is} \alpha_{js} K(\mathbf{X}_i, \mathbf{X}_j).$$

That is $\boldsymbol{\alpha}'_s \mathbf{K} \boldsymbol{\alpha}_s = \lambda_s^2 N^2$. Because of (17.4) this implies that $\boldsymbol{\alpha}'_s \boldsymbol{\alpha}_s = \lambda_s N$.

In order to apply this PCA procedure we need to compute the eigenvectors of matrix \mathbf{K} corresponding to its largest eigenvalues. This will give us vectors $\boldsymbol{\alpha}_s$ and numbers λ_s . For a data point $\mathbf{X} \in \mathbb{R}^m$ its s -PCA component is $\mathbf{e}'_s \mathbf{h}(\mathbf{X})$. By (17.1) we have

$$\mathbf{e}'_s \mathbf{h}(\mathbf{X}) = \frac{1}{\lambda_s N} \sum_{i=1}^N \alpha_{is} \mathbf{h}(\mathbf{X}_i)' \mathbf{h}(\mathbf{X}) = \frac{1}{\lambda_s N} \sum_{i=1}^N \alpha_{is} K(\mathbf{X}_i, \mathbf{X}).$$

When $N^{-1} \sum_{i=1}^N \mathbf{h}(\mathbf{X}_i) \neq \mathbf{0}$ we can make the following correction to the matrix \mathbf{K} :

$$\begin{aligned} \tilde{\mathbf{K}}_{ij} &= [\mathbf{h}(\mathbf{X}_i) - N^{-1} \sum_{k=1}^N \mathbf{h}(\mathbf{X}_k)]' [\mathbf{h}(\mathbf{X}_j) - N^{-1} \sum_{\ell=1}^N \mathbf{h}(\mathbf{X}_\ell)] \\ &= \mathbf{K}_{ij} - N^{-1} \sum_{k=1}^N \mathbf{K}_{ki} - N^{-1} \sum_{k=1}^N \mathbf{K}_{kj} + N^{-2} \sum_{k=1}^N \sum_{\ell=1}^N \mathbf{K}_{k\ell}. \end{aligned}$$

18 Correlation analysis

18.1 Partial correlation

Let X, Y and Z be random variables. Partial correlation between X and Y given Z , denoted $\text{Corr}(X, Y|Z)$ or $\rho_{XY.Z}$, is defined as the correlation between residuals of X and Y regressed on Z . That is, let us consider regression X on Z . Without loss of generality we can assume that $\mathbb{E}[X] = \mathbb{E}[Y] = \mathbb{E}[Z] = 0$. The regression is obtained by solving

$$\min_{\beta} \mathbb{E}[(X - \beta Z)^2].$$

Solution of this problem is $\beta = \text{Cov}(X, Z)/\text{Var}(Z) = \text{Corr}(X, Z)$. Hence the partial correlation is

$$\text{Corr}(X, Y|Z) = \text{Corr}(X - \rho_{XZ}Z, Y - \rho_{YZ}Z) = \frac{\rho_{XY} - \rho_{XZ}\rho_{YZ}}{\sqrt{1 - \rho_{XZ}^2} \sqrt{1 - \rho_{YZ}^2}},$$

where $\rho_{XZ} = \text{Corr}(X, Z)$ and $\rho_{YZ} = \text{Corr}(Y, Z)$.

In similar way partial correlation between random variables X and Y given random variables Z_1, Z_2, \dots, Z_n , is defined. That is, suppose that $\mathbb{E}[X] = \mathbb{E}[Y] = \mathbb{E}[Z_1] = \dots = \mathbb{E}[Z_n] = 0$. Consider the problem

$$\min_{\boldsymbol{\beta}} \mathbb{E}[(X - \boldsymbol{\beta}' \mathbf{Z})^2].$$

Solution of this problem is $\boldsymbol{\beta} = \boldsymbol{\Sigma}_Z^{-1} \boldsymbol{\Sigma}_{ZX}$, where $\boldsymbol{\Sigma}_Z$ is the covariance matrix of random vector $\mathbf{Z} = (Z_1, \dots, Z_n)$ and $\boldsymbol{\Sigma}_{XZ} = \text{Cov}(X, \mathbf{Z})$. Hence

$$\text{Corr}(X, Y|\mathbf{Z}) = \text{Corr}(X - \boldsymbol{\Sigma}_{XZ} \boldsymbol{\Sigma}_Z^{-1} \mathbf{Z}, Y - \boldsymbol{\Sigma}_{YZ} \boldsymbol{\Sigma}_Z^{-1} \mathbf{Z}).$$

18.2 Canonical correlation analysis

Consider random vectors $\mathbf{X} = (X_1, \dots, X_p)'$ and $\mathbf{Y} = (Y_1, \dots, Y_q)'$. Let $\boldsymbol{\mu}_1 = \mathbb{E}[\mathbf{X}]$ and $\boldsymbol{\mu}_2 = \mathbb{E}[\mathbf{Y}]$, and $\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$ be the covariance matrix of $(\mathbf{X}', \mathbf{Y}')$, i.e., $\boldsymbol{\Sigma}_{11} = \text{Cov}(\mathbf{X})$, $\boldsymbol{\Sigma}_{22} = \text{Cov}(\mathbf{Y})$ and $\boldsymbol{\Sigma}_{12} = \text{Cov}(\mathbf{X}, \mathbf{Y})$. Consider random variables $U = \mathbf{a}' \mathbf{X}$ and $V = \mathbf{b}' \mathbf{Y}$ for some vectors $\mathbf{a} \in \mathbb{R}^p$ and $\mathbf{b} \in \mathbb{R}^q$. We want to solve the problem

$$\max_{\mathbf{a}, \mathbf{b}} \text{Corr}(U, V). \quad (18.1)$$

Suppose for the moment that $\Sigma_{11} = \mathbf{I}_p$ and $\Sigma_{22} = \mathbf{I}_q$. Then $\text{Cov}(U, V) = \mathbf{a}'\Sigma_{12}\mathbf{b}$ and $\text{Var}(U) = \mathbf{a}'\mathbf{a}$, $\text{Var}(V) = \mathbf{b}'\mathbf{b}$. Hence problem (18.1) becomes

$$\max_{\mathbf{a}, \mathbf{b}} \frac{\mathbf{a}'\Sigma_{12}\mathbf{b}}{\sqrt{\mathbf{a}'\mathbf{a}}\sqrt{\mathbf{b}'\mathbf{b}}}. \quad (18.2)$$

Note that for a given vector \mathbf{w} , the maximum of $\mathbf{w}'\mathbf{b}$ subject to $\|\mathbf{b}\| = 1$ is attained at $\bar{\mathbf{b}} = \mathbf{w}/\|\mathbf{w}\|$. Therefore for given \mathbf{a} the maximum in (18.2) is attained at $\mathbf{b} = \Sigma_{21}\mathbf{a}$. Hence with respect to \mathbf{a} problem (18.2) becomes

$$\max_{\mathbf{a}} \left\{ \frac{\mathbf{a}'\Sigma_{12}\Sigma_{21}\mathbf{a}}{\sqrt{\mathbf{a}'\mathbf{a}}\sqrt{\mathbf{a}'\Sigma_{12}\Sigma_{21}\mathbf{a}}} = \sqrt{\frac{\mathbf{a}'\Sigma_{12}\Sigma_{21}\mathbf{a}}{\mathbf{a}'\mathbf{a}}} \right\}. \quad (18.3)$$

Optimal solution $\bar{\mathbf{a}}$ of problem (18.3) is given by the eigenvector of matrix $\Sigma_{12}\Sigma_{21}$ corresponding to its largest eigenvalue λ_1 , and the maximum in (18.1) is equal to $\sqrt{\lambda_1}$. Similar the optimal $\bar{\mathbf{b}}$ is given by the eigenvector of matrix $\Sigma_{21}\Sigma_{12}$ corresponding to its largest eigenvalue λ_1 . Note that

$$\Sigma_{21}\Sigma_{12}\Sigma_{21}\bar{\mathbf{a}} = \lambda_1\Sigma_{21}\bar{\mathbf{a}},$$

and hence $\bar{\mathbf{b}} = \Sigma_{21}\bar{\mathbf{a}}$.

In general let $\mathbf{c} = \Sigma_{11}^{1/2}\mathbf{a}$ and $\mathbf{d} = \Sigma_{22}^{1/2}\mathbf{b}$. Then

$$\text{Corr}(U, V) = \frac{\mathbf{a}'\Sigma_{12}\mathbf{b}}{\sqrt{\mathbf{a}'\Sigma_{11}\mathbf{a}}\sqrt{\mathbf{b}'\Sigma_{22}\mathbf{b}}} = \frac{\mathbf{c}'\Sigma_{11}^{-1/2}\Sigma_{12}\Sigma_{22}^{-1/2}\mathbf{d}}{\sqrt{\mathbf{c}'\mathbf{c}}\sqrt{\mathbf{d}'\mathbf{d}}}.$$

Hence the maximum is attained at $\bar{\mathbf{c}}$ given by the eigenvector of $\Sigma_{11}^{-1/2}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1/2}$ corresponding to its largest eigenvalue λ_1 , and at $\bar{\mathbf{d}}$ given by the eigenvector of $\Sigma_{22}^{-1/2}\Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1/2}$ corresponding to its largest eigenvalue λ_1 , and

$$\bar{\mathbf{d}} = \Sigma_{11}^{-1/2}\Sigma_{12}\Sigma_{22}^{-1/2}\bar{\mathbf{c}}.$$

We have that

$$\Sigma_{11}^{-1/2}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1/2}\bar{\mathbf{c}} = \lambda_1\bar{\mathbf{c}}$$

and $\bar{\mathbf{c}} = \Sigma_{11}^{1/2}\bar{\mathbf{a}}$. Hence

$$\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\bar{\mathbf{a}} = \lambda_1\bar{\mathbf{a}}, \quad (18.4)$$

and similarly

$$\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\bar{\mathbf{b}} = \lambda_1\bar{\mathbf{b}}. \quad (18.5)$$

Let $\mathbf{a}_1 = \bar{\mathbf{a}}$ and $\mathbf{b}_1 = \bar{\mathbf{b}}$, and $U_1 = \mathbf{a}'_1\mathbf{X}$ and $V_1 = \mathbf{b}'_1\mathbf{Y}$. At the second stage we want to find $U_2 = \mathbf{a}'_2\mathbf{X}$ and $V_2 = \mathbf{b}'_2\mathbf{Y}$ such that $\text{Cov}(U_2, U_1) = 0$, $\text{Cov}(V_2, V_1) = 0$ and $\text{Corr}(U_2, V_2)$ is maximized. Consider $\mathbf{c}_2 = \Sigma_{11}^{1/2}\mathbf{a}_2$ and $\mathbf{d}_2 = \Sigma_{22}^{1/2}\mathbf{b}_2$. Then

$$\text{Cov}(U_2, U_1) = \mathbf{a}'_2\Sigma_{11}\mathbf{a}_1 = \mathbf{c}'_2\Sigma_{11}^{-1/2}\Sigma_{11}\Sigma_{11}^{-1/2}\mathbf{c}_1 = \mathbf{c}'_2\mathbf{c}_1.$$

Hence $\text{Cov}(U_2, U_1) = 0$ iff $\mathbf{c}'_2\mathbf{c}_1 = 0$. Therefore the second stage problem is

$$\max_{\mathbf{c}} \frac{\mathbf{c}'\Sigma_{11}^{-1/2}\Sigma_{12}\Sigma_{22}^{-1/2}\mathbf{d}}{\sqrt{\mathbf{c}'\mathbf{c}}\sqrt{\mathbf{d}'\mathbf{d}}} \text{ subject to } \mathbf{c}'\mathbf{c}_1 = 0.$$

The maximum is attained at $\bar{\mathbf{c}}$ given by the eigenvector of $\Sigma_{11}^{-1/2}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1/2}$ corresponding to its second largest eigenvalue λ_2 . And so on.

19 Gaussian Mixture Models

Let $y_i \in \{1, \dots, K\}$ be one of K possible labels for data point \mathbf{X}_i , $i = 1, \dots, N$. Assume that the pdf of the data $f(\mathbf{x}_i, y_i) = f(\mathbf{x}_i|y_i)p(y_i)$, is defined as follows: $p(y_i = k) = \pi_k$, $k = 1, \dots, K$, and the conditional distributions $f(\mathbf{x}_i|y_i = k) \sim \mathcal{N}_m(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ are normal. The corresponding log-likelihood function is

$$\ell_n(\boldsymbol{\theta}) = \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k \phi(\mathbf{X}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right),$$

where

$$\phi(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{m/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -(\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) / 2 \right\},$$

and $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K)$ is vector of parameters.

EM (Expectation-Maximization) algorithm

Initialize the means $\boldsymbol{\mu}_k$, covariances $\boldsymbol{\Sigma}_k$ and mixing coefficients π_k , $k = 1, \dots, K$.

The Expectation step (E-step) Given current estimates of the parameters π_1, \dots, π_K , $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$, $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K$, evaluate (by the Bayes rule) the corresponding posterior probabilities of data point \mathbf{X}_i being in cluster $k \in \{1, \dots, K\}$:

$$w_{ik} = \frac{\pi_k \phi(\mathbf{X}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \phi(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}, \quad i = 1, \dots, N.$$

Note that $\sum_{k=1}^K w_{ik} = 1$ for all i .

The Maximization step (M-step) For $k = 1, \dots, K$, set $N_k = \sum_{i=1}^N w_{ik}$, and update $\pi_k^{new} = N_k/N$, $\boldsymbol{\mu}_k^{new} = N_k^{-1} \sum_{i=1}^N w_{ik} \mathbf{X}_i$, and

$$\boldsymbol{\Sigma}_k^{new} = N_k^{-1} \sum_{i=1}^N w_{ik} (\mathbf{X}_i - \boldsymbol{\mu}_k) (\mathbf{X}_i - \boldsymbol{\mu}_k)'.$$

Note that $\sum_{k=1}^K N_k = \sum_{i=1}^N \sum_{k=1}^K w_{ik} = N$.

20 Von Mises statistical functionals

Let $\mathbf{X}_1, \dots, \mathbf{X}_N$ be an iid sample of random vectors with probability distribution (probability measure) $\mathbf{X}_i \sim F(\cdot)$. With the sample is associated the so called empirical probability measure (distribution) $\hat{F}_N = N^{-1} \sum_{i=1}^N \delta_{\mathbf{X}_i}$, where δ_x denotes probability measure of mass 1 at the point x . When X_1, \dots, X_N are random numbers, the empirical cdf $\hat{F}_N(x) = \frac{\#(X_i \leq x)}{N}$. That is, if the sample is arranged in the increasing order $X_{(1)} \leq \dots \leq X_{(N)}$, then $\hat{F}_N(x) = 0$ for $x < X_{(1)}$, $\hat{F}_N(x) = 1/N$ for $X_{(1)} \leq x < X_{(2)}$, $\hat{F}_N(x) = 2/n$ for $X_{(2)} \leq x < X_{(3)}$, and so on.

Function $\theta = T(F)$ of the distribution F is called statistical functional. Its sample estimate is $\hat{\theta} = T(\hat{F}_N)$. Consider the following examples.

- Expectation of a function:

$$T(F) = \mathbb{E}_F[h(\mathbf{X})] = \int h(\mathbf{x}) dF(\mathbf{x}).$$

Its sample estimate

$$T(\hat{F}_N) = \mathbb{E}_{\hat{F}_N}[h(\mathbf{X})] = N^{-1} \sum_{i=1}^N h(\mathbf{X}_i).$$

- Variance

$$T(F) = \text{Var}(X) = \mathbb{E}_F[X^2] - (\mathbb{E}_F[X])^2.$$

Its sample estimate

$$T(\hat{F}_N) = N^{-1} \sum_{i=1}^N X_i^2 - \bar{X}^2 = N^{-1} \sum_{i=1}^N (X_i - \bar{X})^2.$$

- Median⁸ it is defined

$$T(F) = F^{-1}(1/2).$$

Its sample estimate

$$T(\hat{F}_N) = \hat{F}_N^{-1}(1/2).$$

- Solution of equation $\mathbb{E}_F[g(\mathbf{X}, \boldsymbol{\theta})] = 0$. Its sample estimate is obtained as solution of equation $\mathbb{E}_{\hat{F}_N}[g(\mathbf{X}, \hat{\boldsymbol{\theta}})] = 0$, which is $\sum_{i=1}^N g(\mathbf{X}_i, \hat{\boldsymbol{\theta}}) = 0$.

It is known (Glivenko-Cantelli Theorem) that the empirical cdf $\hat{F}_N(x)$ converges w.p.1 to $F(x)$ uniformly in $x \in \mathbb{R}$, that is $\sup_{x \in \mathbb{R}} |\hat{F}_N(x) - F(x)|$ converges w.p.1 to 0 as N tends to infinity. If $T(\cdot)$ is continuous (in a certain sense), it follows then that $T(\hat{F}_N)$ converges to $T(F)$ w.p.1, i.e., $\hat{\theta} = T(\hat{F}_N)$ is a consistent estimator of $\theta = T(F)$.

Asymptotic normality Consider probability distributions F and G . Their convex combination is

$$(1-t)F + tG = F + t(G - F), \quad t \in [0, 1].$$

The directional derivative of $T(\cdot)$ at F in the direction $G - F$ is

$$T'(F, G - F) = \lim_{t \downarrow 0} \frac{T(F + t(G - F)) - T(F)}{t}.$$

That is, $T'(F, G - F)$ is the right side derivative of $F_t := (1-t)F + tG$ at $t = 0$. Let $G = \hat{F}_N$. Then

$$T'(F, \hat{F}_N - F) = T' \left(F, N^{-1} \sum_{i=1}^N \delta_{X_i} - F \right) = T' \left(F, N^{-1} \sum_{i=1}^N [\delta_{X_i} - F] \right).$$

Suppose further that $T'(F, \cdot)$ is linear (as a function of the direction), then it follows by the above that

$$T'(F, \hat{F}_N - F) = N^{-1} \sum_{i=1}^N T'(F, \delta_{X_i} - F).$$

Now we use the following approximation

$$\hat{\theta} - \theta = T(\hat{F}_N) - T(F) \approx T'(F, \hat{F}_N - F) = N^{-1} \sum_{i=1}^N IC_{T,F}(\mathbf{X}_i),$$

⁸Quantile $F^{-1}(\alpha)$, $\alpha \in (0, 1)$, is defined by the equation $F(x) = \alpha$. Solution of this equation could be not unique or does not exist if the cdf $F(\cdot)$ is discontinuous. Therefore the left side quantile is defined as $\inf\{x : F(x) \geq \alpha\}$, and the right side quantile is defined as $\sup\{x : F(x) \leq \alpha\}$. If the left side and right side quantiles are different from each other, sometimes their average is taken as the corresponding quantile.

where

$$IC_{T,F}(\mathbf{x}) := T'(F, \delta_x - F) = \lim_{t \downarrow 0} \frac{T((1-t)F + t\delta_x) - T(F)}{t},$$

is the so called Influence Curve (or Influence Function).

Let us note that $\mathbb{E}_F [IC_{T,F}(\mathbf{X})] = 0$. Indeed suppose for the moment that F has discrete distribution, i.e., $F = \sum_{i=1}^m p_i \delta_{x_i}$ for some \mathbf{x}_i and probabilities $p_i > 0$. Then

$$\mathbb{E}_F [IC_{T,F}(\mathbf{X})] = \sum_{i=1}^m p_i IC_{T,F}(\mathbf{x}_i) = \sum_{i=1}^m p_i T'(F, \delta_{x_i} - F) = T'\left(F, \sum_{i=1}^m p_i \delta_{x_i} - F\right),$$

where the last equality holds by linearity of $T'(F, \cdot)$ and since $\sum_{i=1}^m p_i = 1$. Since $\sum_{i=1}^m p_i \delta_{x_i} = F$ and $T'(F, F - F) = 0$, it follows that $\mathbb{E}_F [IC_{T,F}(\mathbf{X})] = 0$.

By the above

$$N^{1/2} \left[T(\hat{F}_N) - T(F) \right] \approx N^{-1/2} \sum_{i=1}^N IC_{T,F}(\mathbf{X}_i). \quad (20.1)$$

Since $\mathbb{E}_F [IC_{T,F}(\mathbf{X}_i)] = 0$, we have by the CLT that $N^{-1/2} \sum_{i=1}^N IC_{T,F}(\mathbf{X}_i)$ converges in distribution to normal with zero mean and variance

$$\sigma_{T,F}^2 = \mathbb{E}_F [IC_{T,F}(\mathbf{X})^2] = \text{Var}_F [IC_{T,F}(\mathbf{X})].$$

This suggests that $N^{1/2} [T(\hat{F}_N) - T(F)]$ converges in distribution to normal $\mathcal{N}(0, \sigma^2)$ with $\sigma^2 = \text{Var}_F [IC_{T,F}(\mathbf{X})]$.

These derivations of asymptotic normality of $T(\hat{F}_N)$ are somewhat heuristic since the approximation (20.1) is not rigorously justified. Nevertheless it usually gives correct results, which could be proved by ad hoc methods, and is routinely used in applications.

For example, consider the median functional $T(F) = F^{-1}(1/2)$ (here F is the cumulative distribution function). Suppose that the (population) median $\mathbf{m} = F^{-1}(1/2)$ is uniquely defined and the distribution has density $dF(\mathbf{m})/dx = f(\mathbf{m})$ at $x = \mathbf{m}$.

Let us compute the directional derivative $T'(F, G - F)$ for some cdf G . Let $F_t = (1-t)F + tG$ and consider $T(F_t) = F_t^{-1}(1/2)$. We have that $F_t(T(F_t)) = 1/2$, i.e.,

$$(1-t)F(T(F_t)) + tG(T(F_t)) = 1/2.$$

Computing derivative of the above with respect to t gives

$$-F(T(F_t)) + (1-t) \frac{dF(T(F_t))}{dt} + G(T(F_t)) + t \frac{dG(T(F_t))}{dt} = 0. \quad (20.2)$$

At $t = 0$ we have that $F_0 = F$ and

$$\left. \frac{dF(T(F_t))}{dt} \right|_{t=0} = f(\mathbf{m}) \left. \frac{dT(F_t)}{dt} \right|_{t=0}. \quad (20.3)$$

Equation (20.2) (for $t = 0$) together with (20.2) imply that

$$-F(\mathbf{m}) + G(\mathbf{m}) + f(\mathbf{m}) \left. \frac{dT(F_t)}{dt} \right|_{t=0} = 0,$$

and hence (since $F(\mathbf{m}) = 1/2$)

$$T'(F, G - F) = \left. \frac{dT(F_t)}{dt} \right|_{t=0} = \frac{1/2 - G(\mathbf{m})}{f(\mathbf{m})}$$

We obtain that

$$IC_{T,F}(x) = T'(F, \delta_x - F) = \frac{1/2 - \delta_x(\mathbf{m})}{f(\mathbf{m})},$$

where δ_x is the cdf such that $\delta_x(t) = 0$ for $t < x$, and $\delta_x(t) = 1$ for $t \geq x$.

Note that $\mathbb{E}_F [IC_{T,F}(X)] = 0$ (as it should be), since $\mathbb{E}_F [\delta_X(\mathbf{m})] = P(X \leq \mathbf{m}) = 1/2$. Also $\text{Var}_F [\delta_X(\mathbf{m})] = 1/2 - 1/4 = 1/4$ and hence

$$\text{Var}_F [IC_{T,F}(X)] = \frac{\text{Var}_F [\delta_X(\mathbf{m})]}{f(\mathbf{m})^2} = \frac{1}{4f(\mathbf{m})^2}.$$

We obtain that $N^{1/2}[T(\hat{F}_N) - T(F)]$ converges in distribution to normal $\mathcal{N}(0, \frac{1}{4f(\mathbf{m})^2})$. That is the sample median has approximately normal distribution with variance $\frac{1}{4Nf(\mathbf{m})^2}$, provided that the population median \mathbf{m} is defined uniquely and the distribution has density $f(\mathbf{m}) = dF(\mathbf{m})/dx$ at $x = \mathbf{m}$.

For example suppose that variables X_i have normal distribution $\mathcal{N}(\mu, \sigma^2)$. In that case the median $\mathbf{m} = \mu$. Asymptotic variance of the sample median is $N^{-1}\sigma^2(\pi/2)$, while variance of \bar{X} is $N^{-1}\sigma^2$. In that case \bar{X} is a better estimator of $\mathbf{m} = \mu$.

However, suppose now that X_i have Laplace distribution with $f(x, \theta) = \frac{1}{2}e^{-|x-\theta|}$, $\theta \in \mathbb{R}$. Then θ is the mean and median of the distribution, and $N^{1/2}(\hat{\theta} - \theta)$ converges in distribution to normal $\mathcal{N}(0, 1)$. We have here that $\text{Var}(X_i) = 2$ and hence variance of \bar{X} is $2N^{-1}$, while the asymptotic variance of the sample median is N^{-1} . It is also interesting to note that the MLE $\hat{\theta}$ is the sample median. Now $\partial \log f(x, \theta)/\partial \theta$ is equal 1 if $\theta < x$ and -1 if $\theta > x$. Thus $\partial^2 \log f(x, \theta)/\partial \theta^2 = 0$ for $\theta \neq x$, and $\partial^2 \log f(x, \theta)/\partial \theta^2$ is not defined for $\theta = x$. Hence formula (8.4) for the information number cannot be applied, i.e., the situation here is not standard.

As another example suppose that Y has Cauchy distribution, i.e., $Y = V/W$ with independent $V \sim \mathcal{N}(0, 1)$ and $W \sim \mathcal{N}(0, 1)$. Cauchy distribution has pdf $f_Y(y) = \frac{1}{\pi(1+y^2)}$. Therefore in that case asymptotic variance of the sample median is $N^{-1}\pi^2/4$. On the other hand, $\mathbb{E}|Y| = +\infty$ and the average \bar{X} has the same distribution as $\mathbf{m} + Y$ for any sample size N , and will not converge to \mathbf{m} as $N \rightarrow \infty$.

Finite sample interpretation of the influence curve. By adding one more observation X_{N+1} to sample X_1, \dots, X_N , we have that

$$\hat{F}_{N+1}(\cdot) = \frac{N}{N+1}\hat{F}_N(\cdot) + \frac{1}{N+1}\delta_{X_{N+1}}(\cdot) = (1-t)\hat{F}_N(\cdot) + t\delta_{X_{N+1}}(\cdot),$$

where $t = 1/(N+1)$. Hence we can write

$$\hat{\theta}_{N+1} \approx \hat{\theta}_N + \frac{1}{N+1}IC_{T, \hat{F}_N}(X_{N+1}).$$

This shows sensitivity of the estimator to one observation. If $\text{Var}_F [IC_{T,F}(X)]$ is large, the estimator $T(\hat{F}_N)$ can be sensitive just to one observation.

21 Bootstrap

21.1 Jackknife bias estimation.

Consider an estimator $\hat{\theta} = \hat{\theta}(\mathbf{X}_1, \dots, \mathbf{X}_N)$. Denote

$$\hat{\theta}_{-i} = \hat{\theta}(\mathbf{X}_1, \dots, \mathbf{X}_{i-1}, \mathbf{X}_{i+1}, \dots, \mathbf{X}_N), \quad i = 1, \dots, N,$$

i.e., $\hat{\theta}_{-i}$ is obtained by removing data point \mathbf{X}_i from calculation of $\hat{\theta}$. Let $\bar{\theta} = N^{-1} \sum_{i=1}^N \hat{\theta}_{-i}$.

The Jackknife estimator of the bias is $(N-1)(\bar{\theta} - \hat{\theta})$. Bias corrected version Jackknife estimator

$$\hat{\theta}_{jack} = \hat{\theta} - (N-1)(\bar{\theta} - \hat{\theta}) = N\hat{\theta} - (N-1)(\bar{\theta}). \quad (21.1)$$

Theoretical justification. Suppose that

$$\mathbb{E}_{\theta}[\hat{\theta}] = \theta + N^{-1}a(\theta),$$

i.e., bias $b_{\theta}(\hat{\theta}) = N^{-1}a(\theta)$, of $\hat{\theta}$, is of order $O(1/N)$. Then $\mathbb{E}[\hat{\theta}_{-i}] = \theta + (N-1)^{-1}a(\theta)$ and hence

$$\mathbb{E}[\bar{\theta}] = N^{-1} \sum_{i=1}^n \mathbb{E}[\hat{\theta}_{-i}] = \theta + (N-1)^{-1}a(\theta),$$

and thus

$$\mathbb{E}[\hat{\theta} - \bar{\theta}] = N^{-1}a(\theta) - (N-1)^{-1}a(\theta) = [N(N-1)]^{-1}a(\theta).$$

It follows that

$$\mathbb{E}[(N-1)(\hat{\theta} - \bar{\theta})] = -N^{-1}a(\theta),$$

and hence $\mathbb{E}[\hat{\theta}_{jack}] = \theta$, i.e. $\hat{\theta}_{jack}$ is an unbiased estimator of θ .

21.2 Bootstrap method

The idea of resampling used in the Jackknife estimation is further extended in the Bootstrap method. Let $\hat{\theta} = \hat{\theta}(\mathbf{X}_1, \dots, \mathbf{X}_N)$ be an estimator which is a function of sample $\mathbf{X}_1, \dots, \mathbf{X}_N$. Suppose that we want to evaluate statistical properties of that estimator without assuming a parametric model. For example we would like to construct two sided 95% confidence interval for this estimator. This means that we need to evaluate 2.5% and 97.5% quantiles of the distribution of $\hat{\theta}$. Note that both quantiles are functions of the true distribution F of the sample. If we knew the true distribution F we can proceed by using the so called Monte Carlo sampling techniques. That is, we generate a sample $\tilde{\mathbf{X}}_1, \dots, \tilde{\mathbf{X}}_N$ from F and compute $\tilde{\theta} = \hat{\theta}(\tilde{\mathbf{X}}_1, \dots, \tilde{\mathbf{X}}_N)$. We repeat this procedure independently M times. In that way we generate M independent replications $\tilde{\theta}_1, \dots, \tilde{\theta}_M$ of the random variable θ . Consequently for sufficiently large M , we can accurately reconstruct the true distribution of $\hat{\theta}$, and hence to evaluate the required quantiles, or some other parameters. For example we can estimate variance of $\hat{\theta}$ as

$$\widehat{\text{Var}}(\hat{\theta}) = \frac{1}{M-1} \sum_{m=1}^M (\tilde{\theta}_m - \bar{\tilde{\theta}})^2,$$

where $\bar{\tilde{\theta}} = \frac{1}{M} \sum_{m=1}^M \tilde{\theta}_m$.

Of course the true distribution F is not known. So we replace it by the empirical distribution $\hat{F}_N = N^{-1} \sum_{i=1}^N \delta_{X_i}$. Then we proceed by generating a random sample $\mathbf{X}_1^*, \dots, \mathbf{X}_N^*$ from \hat{F}_N and compute $\hat{\theta}^* = \hat{\theta}(\mathbf{X}_1^*, \dots, \mathbf{X}_N^*)$. We repeat this procedure M times to obtain values $\hat{\theta}_1^*, \dots, \hat{\theta}_M^*$,

which can be used to estimate quantity of interest. Generating a sample $\mathbf{X}_1^*, \dots, \mathbf{X}_N^*$ from \hat{F}_N means resampling from the data set (sample) $\mathbf{X}_1, \dots, \mathbf{X}_N$. That is, an element \mathbf{X}_i^* is chosen at random from the set $\{\mathbf{X}_1, \dots, \mathbf{X}_N\}$. This is repeated N times with replacement, to generate one realization $\mathbf{X}_1^*, \dots, \mathbf{X}_N^*$. So each element of the generated sample $\mathbf{X}_1^*, \dots, \mathbf{X}_N^*$ coincides with some element of the original sample $\mathbf{X}_1, \dots, \mathbf{X}_N$.

This procedure is easy to implement and does not require any modelling assumptions. On the other hand, it is solely based on the sample (the data) $\mathbf{X}_1, \dots, \mathbf{X}_N$ and can be very sensitive to outliers. Its theoretical analysis is quite sophisticated and is based of theory of statistical estimators of functionals $\theta = T(F)$.

22 Robust statistics

Let $\rho : \mathbb{R} \rightarrow \mathbb{R}_+$ be a convex function such that $\rho(0) = 0$. Consider the problem

$$\min_{\theta} \mathbb{E}[\rho(X - \theta)].$$

For example if $\rho(t) = t^2$ this becomes the least squares problem, its solution is $\theta^* = \mathbb{E}[X]$. Another example $\rho(t) = |t|$. In that case solution θ^* is the median of the distribution of X . The sample median is much less sensitive to outliers than the average. Another example which tries to combine local efficiency of least squares with robustness of absolute value deviations is $\rho(t) = t^2$ for $|t| \leq 2$, and $\rho(t) = |t| + 2$ for $|t| \geq 2$. So when observation in the interval $[-2, 2]$ it works like least squares, outside that interval it could deal with outliers as the absolute deviation method.

As another example for $\alpha \in (0, 1)$ let

$$\rho_{\alpha}(t) = \begin{cases} -(1 - \alpha)t & \text{if } t \leq 0, \\ \alpha t & \text{if } t \geq 0. \end{cases} \quad (22.1)$$

Consider the problem

$$\min_{\theta \in \mathbb{R}} \mathbb{E}[\rho_{\alpha}(X - \theta)]. \quad (22.2)$$

We have that

$$\partial \mathbb{E}[\rho_{\alpha}(X - \theta)] / \partial \theta = \mathbb{E}[\partial \rho_{\alpha}(X - \theta) / \partial \theta],$$

with $\partial \rho_{\alpha}(X - \theta) / \partial \theta$ is equal to $(1 - \alpha)$ for $X - \theta < 0$, and $-\alpha$ for $X - \theta > 0$. Suppose that $F(x)$ is continuous at $x = \theta$. It follows that

$$\partial \mathbb{E}[\rho_{\alpha}(X - \theta)] / \partial \theta = F(\theta) - \alpha,$$

where $F(x) = \mathbf{Prob}(X \leq x)$ is the cdf of X . Thus the quantile $\theta = F^{-1}(\alpha)$ is the optimal solution of problem (22.2). In particular for $\alpha = 1/2$, $\rho_{\alpha}(t) = \frac{1}{2}|t|$ and solution of problem (22.2) is the median of the distribution. As it was discussed in section 20, the left side $\inf\{x : F(x) \geq \alpha\}$, and the right side $\sup\{x : F(x) \leq \alpha\}$ quantiles can be different from each other. In that case optimal solution of problem (22.2) can be any point between the left side and right side quantiles.

22.1 Quantile regression

Quantile regression use function $\rho_{\alpha}(\cdot)$, defined in (22.1), to fit linear model to the data. That is, consider the problem

$$\min_{\beta \in \mathbb{R}^{k+1}} \mathbb{E}[\rho_{\alpha}(Y - \beta' \mathbf{X})], \quad (22.3)$$

where Y and $\mathbf{X} = (1, X_1, \dots, X_k)'$ are random variables. Given data Y_i and $\mathbf{X}_i = (1, X_{i1}, \dots, X_{ik})'$, $i = 1, \dots, N$, the sample counterpart of problem (22.3) is

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^{k+1}} \sum_{i=1}^N \rho_\alpha(Y_i - \boldsymbol{\beta}' \mathbf{X}_i), \quad (22.4)$$

The solution $\hat{\boldsymbol{\beta}}$ of problem (22.4) can be viewed as an estimator of the solution of problem (22.3). For $\alpha = 1/2$ this becomes the least absolute deviations method for solving linear regression. Both problems (22.3) and (22.4) could have more than one optimal solution.

Problem (22.4) can be written as the linear program

$$\begin{aligned} \min_{\boldsymbol{\beta}, \mathbf{v}^+, \mathbf{v}^-} \quad & \sum_{i=1}^N (1 - \alpha)v_i^- + \alpha v_i^+ \\ \text{s.t.} \quad & Y_i - \boldsymbol{\beta}' \mathbf{X}_i = v_i^+ - v_i^-, \quad i = 1, \dots, N, \\ & v_i^- \geq 0, v_i^+ \geq 0, \quad i = 1, \dots, N. \end{aligned}$$

It is possible to show that under some regularity conditions, in particular if Y has pdf $f_Y(\cdot)$, $\sqrt{N}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$ converges in distribution to normal with zero mean vector and covariance matrix $\alpha(1 - \alpha)\boldsymbol{\Psi}^{-1}\boldsymbol{\Omega}\boldsymbol{\Psi}^{-1}$ with $\boldsymbol{\Psi} = \mathbb{E}[f_Y(\boldsymbol{\beta}' \mathbf{X})\mathbf{X}\mathbf{X}']$ and $\boldsymbol{\Omega} = \mathbb{E}[\mathbf{X}\mathbf{X}']$.

23 Bayes estimators

Recall Bayes' formula: if $\{A_i\}$ is a partition of the sample space and B is an event such that $P(B) \neq 0$ then

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_j P(B|A_j)P(A_j)}.$$

Let $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_N)$ be a sample with $\mathbf{X} \sim f(\mathbf{x}, \boldsymbol{\theta})$. Suppose that $\boldsymbol{\theta}$ is random with pdf $\pi(\boldsymbol{\theta})$, referred to as the *prior distribution*. Denote by $f(\mathbf{x}|\boldsymbol{\theta})$ the sampling distribution conditional on $\boldsymbol{\theta}$. Then the joint distribution of \mathbf{X} and $\boldsymbol{\theta}$ is $f(\mathbf{x}, \boldsymbol{\theta}) = f(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})$. By Bayes' formula, the distribution of $\boldsymbol{\theta}$, conditional on $\mathbf{X} = \mathbf{x}$, is

$$\pi(\boldsymbol{\theta}|\mathbf{x}) = \frac{f(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\int f(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\boldsymbol{\theta}}, \quad (23.1)$$

that is $\pi(\boldsymbol{\theta}|\mathbf{x})$ is proportional to $f(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})$, written $\pi(\boldsymbol{\theta}|\mathbf{x}) \propto f(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})$. The distribution (pdf) $\pi(\boldsymbol{\theta}|\mathbf{x})$ is called the *posterior distribution*.

Example 23.1 Suppose that $X_i \sim \mathcal{N}(\theta, \sigma^2)$ and $\theta \sim \mathcal{N}(\mu, \tau^2)$, where σ^2 , μ and τ^2 are supposed to be known. We have that

$$f(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta}) = \frac{1}{(\sqrt{2\pi}\sigma)^N} e^{-\sum_{i=1}^N (x_i - \theta)^2 / 2\sigma^2} \frac{1}{\sqrt{2\pi}\tau} e^{-(\theta - \mu)^2 / 2\tau^2},$$

and hence

$$f(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta}) \propto \exp \left\{ \frac{-\left(\theta - \frac{\tau^2 \sum_{i=1}^N x_i + \mu\sigma^2}{N\tau^2 + \sigma^2}\right)^2}{2\left(\frac{\sigma^2\tau^2}{N\tau^2 + \sigma^2}\right)} \right\}.$$

It follows that the posterior distribution $\pi(\boldsymbol{\theta}|\mathbf{x})$ is normal with conditional mean

$$\mathbb{E}[\theta|\mathbf{x}] = \frac{\tau^2}{\tau^2 + \sigma^2/N} \bar{x} + \frac{\sigma^2/N}{\tau^2 + \sigma^2/N} \mu$$

and conditional variance

$$\text{Var}(\theta|\mathbf{x}) = \frac{(\sigma^2/N)\tau^2}{\sigma^2/N + \tau^2}.$$

Note that $\mathbb{E}[\theta|\mathbf{x}] - \bar{x}$ tends to 0 and $\text{Var}(\theta|\mathbf{x})$ tends to 0 as $N \rightarrow \infty$. That is, if we view the ‘true’ distribution of the sample as normal with mean θ^* and variance σ^2 , then the average \bar{X} converges in probability to θ^* , and the Bayes estimator converges in probability to θ^* , i.e. for any $\varepsilon > 0$ the probability $\text{Prob}(|\pi(\theta|\mathbf{x}) - \theta^*| > \varepsilon)$ converges to 0 w.p.1 as $N \rightarrow \infty$. The probability is with respect to the true distribution of X_i and the convergence w.p.1 is with respect to the true distribution.

This is a general property of Bayes estimators. If we assume that the true distribution of the sample X_i , $i = 1, \dots$, is $f(x, \theta^*)$ for some $\theta^* \in \Theta$, then (under some regularity conditions) Bayes estimator converges in probability to θ^* for almost every (with respect to the true distribution) sequence X_1, \dots .

In general it may be not easy to compute the posterior distribution. The problem is in calculation of the integral in the right hand side of (23.1). In the above example it was possible to compute the posterior distribution in a closed form, and the posterior distribution was in the same family of normal distributions. Such families of distributions are called *conjugate families*.

Example 23.2 One parameter exponential family

$$f(x|\theta) = \exp[\eta(\theta)T(x) - A(\theta)]h(x),$$

with prior

$$\pi(\theta) \propto \exp[\alpha\eta(\theta) - \beta A(\theta)].$$

Then posterior distribution

$$f(\theta|\mathbf{x})\pi(\theta) \propto \exp[\eta(\theta)(T(\mathbf{x}) + \alpha) - (\beta + 1)A(\theta)]$$

is in the same family of one parameter exponential distributions. □

23.1 Bayesian decisions

Consider a loss function $L(\theta, a)$ (see definition 8.5) and let $\delta(\mathbf{X})$ be a decision rule, e.g., $\delta(\mathbf{X})$ is an estimator of parameter θ . The corresponding risk function is

$$R(\theta, \delta) = \mathbb{E}_\theta[L(\theta, \delta(\mathbf{X}))] = \int L(\theta, \delta(\mathbf{x})) f(\mathbf{x}, \theta) d\mathbf{x}.$$

For example, let $L(\theta, a) = (\theta - a)^2$. Then

$$R(\theta, \delta) = \mathbb{E}_\theta[(\theta - \delta(\mathbf{X}))^2] = \text{Var}_\theta(\delta(\mathbf{X})) + \underbrace{(\mathbb{E}_\theta[\delta(\mathbf{X})] - \theta)^2}_{\text{bias}(\delta(\mathbf{X}))}$$

is the Mean Square Error of the estimator $\delta(\mathbf{X})$ of θ .

Bayes risk, with prior $\pi(\theta)$:

$$B(\pi, \delta) = \mathbb{E}_\pi[R(\theta, \delta)] = \int R(\theta, \delta)\pi(\theta)d\theta = \int \left(\int L(\theta, \delta(\mathbf{x}))f(\mathbf{x}|\theta)d\mathbf{x} \right)\pi(\theta)d\theta.$$

The Bayes rule with respect to the prior $\pi(\theta)$ is

$$\delta^\pi \in \arg \min_{\delta \in \mathcal{D}} B(\pi, \delta), \tag{23.2}$$

where \mathcal{D} is a family of decision rules.

Theorem 23.1 Define

$$r(\mathbf{x}, a) = \int L(\boldsymbol{\theta}, a)\pi(\boldsymbol{\theta}|\mathbf{x})d\boldsymbol{\theta},$$

and let $\delta^\pi(\mathbf{x})$ be a minimizer of $r(\mathbf{x}, a)$, i.e., $\delta^\pi(\mathbf{x}) \in \arg \min_a r(\mathbf{x}, a)$. Suppose that $\delta^\pi \in \mathcal{D}$. Then δ^π is the Bayes rule with respect to π .

Proof. Denote $\mathbf{m}(\mathbf{x}) := \int f(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\boldsymbol{\theta}$. We can write

$$\begin{aligned} B(\pi, \delta) &= \int \left(\int L(\boldsymbol{\theta}, \delta(\mathbf{x}))f(\mathbf{x}|\boldsymbol{\theta})d\mathbf{x} \right) \pi(\boldsymbol{\theta})d\boldsymbol{\theta} \\ &= \int \int L(\boldsymbol{\theta}, \delta(\mathbf{x}))f(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\mathbf{x}d\boldsymbol{\theta} \\ &= \int \int L(\boldsymbol{\theta}, \delta(\mathbf{x}))\pi(\boldsymbol{\theta}|\mathbf{x})\mathbf{m}(\mathbf{x})d\boldsymbol{\theta}d\mathbf{x} \\ &= \int \int r(\mathbf{x}, \delta(\mathbf{x}))\mathbf{m}(\mathbf{x})d\mathbf{x}. \end{aligned}$$

Since $\delta^\pi(\mathbf{x}) \in \arg \min_a r(\mathbf{x}, a)$, we have that for any $\delta \in \mathcal{D}$,

$$r(\mathbf{x}, \delta^\pi(\mathbf{x}))\mathbf{m}(\mathbf{x}) \leq r(\mathbf{x}, \delta(\mathbf{x}))\mathbf{m}(\mathbf{x}).$$

It follows that δ^π is a minimizer of $B(\pi, \delta)$ over $\delta \in \mathcal{D}$, provided $\delta^\pi \in \mathcal{D}$. □

For squared error loss $L(\theta, a) = (\theta - a)^2$, we have

$$r(\mathbf{x}, a) = \int (\theta - a)^2 \pi(\theta|\mathbf{x})d\theta,$$

and hence

$$\delta^\pi(\mathbf{x}) = \int \theta \pi(\theta|\mathbf{x})d\theta = \mathbb{E}_\pi[\theta|\mathbf{X} = \mathbf{x}]$$

is the posterior mean. For the absolute error loss $L(\theta, a) = |\theta - a|$, the Bayes rule $\delta^\pi(\mathbf{x})$ is the posterior median.

Definition 23.1 It is said that decision rule δ' is as good as δ if $R(\boldsymbol{\theta}, \delta') \leq R(\boldsymbol{\theta}, \delta)$ for all $\boldsymbol{\theta} \in \Theta$. Moreover, if $R(\boldsymbol{\theta}, \delta') < R(\boldsymbol{\theta}, \delta)$ for some $\boldsymbol{\theta} \in \Theta$, it is said that δ' is better than δ .

A decision rule $\delta \in \mathcal{D}$ is admissible if there is no $\delta' \in \mathcal{D}$ that is better than δ .

The following example shows that an admissible decision rule can be quite awkward.

Example 23.3 Suppose that X has Binomial distribution $Bin(n, \theta)$, i.e., for $x = 0, 1, \dots, n$,

$$P(X = x) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}, \quad \theta \in (0, 1).$$

Let $\delta(x) = c$ for some constant $c \in (0, 1)$ and all $x = 0, \dots, n$, and $L(\theta, a) = (\theta - a)^2$ be the loss function. Then

$$R(c, \delta) = \sum_{x=0}^n (\delta(x) - c)^2 P(X = x|\theta = c) = 0.$$

Let δ' be as good as δ . Then

$$0 \leq \sum_{x=0}^n (\delta'(x) - c)^2 P(X = x|\theta = c) = R(c, \delta') \leq R(c, \delta) = 0.$$

Therefore $\delta'(x) = c$ for all $x = 0, \dots, n$, and hence $\delta' = \delta$. It follows that δ is admissible. □

Theorem 23.2 Suppose that $R(\boldsymbol{\theta}, \delta)$ is continuous in $\boldsymbol{\theta}$ and for every $\boldsymbol{\theta} \in \Theta$ there is $\varepsilon > 0$ such that $\int_{V_{\varepsilon, \boldsymbol{\theta}}} \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} > 0$, where $V_{\varepsilon, \boldsymbol{\theta}} = \{\boldsymbol{\theta}' \in \Theta : \|\boldsymbol{\theta}' - \boldsymbol{\theta}\| \leq \varepsilon\}$ is the ε -neighborhood of $\boldsymbol{\theta}$. Then δ^π is an admissible decision rule.

Proof. We argue by a contradiction. Suppose that $\delta \in \mathcal{D}$ is a decision rule which is better than δ^π . Since

$$\delta^\pi \in \arg \min_{\delta \in \mathcal{D}} B(\pi, \delta) = \arg \min_{\delta \in \mathcal{D}} \int R(\boldsymbol{\theta}, \delta) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta},$$

we have that

$$\begin{aligned} 0 &\geq \int R(\boldsymbol{\theta}, \delta^\pi) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} - \int R(\boldsymbol{\theta}, \delta) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ &= \int [R(\boldsymbol{\theta}, \delta^\pi) - R(\boldsymbol{\theta}, \delta)] \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}. \end{aligned}$$

On the other hand, since δ is better than δ^π ,

$$R(\boldsymbol{\theta}, \delta^\pi) - R(\boldsymbol{\theta}, \delta) \geq 0, \quad \forall \boldsymbol{\theta} \in \Theta,$$

and there is a point $\boldsymbol{\theta}^* \in \Theta$ such that $R(\boldsymbol{\theta}^*, \delta^\pi) - R(\boldsymbol{\theta}^*, \delta) > 0$. Since $R(\boldsymbol{\theta}, \delta)$ is continuous in $\boldsymbol{\theta}$, there is a neighborhood Ξ of $\boldsymbol{\theta}^*$ and $\gamma > 0$ such that $R(\boldsymbol{\theta}, \delta^\pi) - R(\boldsymbol{\theta}, \delta) \geq \gamma$ for all $\boldsymbol{\theta} \in \Xi$. By the assumption of the theorem there is ε -neighborhood V of $\boldsymbol{\theta}^*$ such that $V \subset \Xi$ and $\int_V \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} > 0$. It follows that

$$\int [R(\boldsymbol{\theta}, \delta^\pi) - R(\boldsymbol{\theta}, \delta)] \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} \geq \gamma \int_V \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} > 0.$$

This gives the required contradiction. \square

Let $T(\mathbf{X})$ be a sufficient statistic for $\boldsymbol{\theta}$, and let $L(\boldsymbol{\theta}, a)$ be a loss function. Suppose that $L(\boldsymbol{\theta}, a)$ is convex in a for all $\boldsymbol{\theta}$. Consider $\delta^*(t) = \mathbb{E}[\delta(\mathbf{X})|T = t]$. Note that by sufficiency of T , $\delta^*(t)$ does not depend on $\boldsymbol{\theta}$. Since $L(\boldsymbol{\theta}, a)$ is convex in a , we have by Jensen inequality

$$\mathbb{E}[L(\boldsymbol{\theta}, \delta(\mathbf{X})|T] \geq L(\boldsymbol{\theta}, \mathbb{E}[\delta(\mathbf{X})|T]) = L(\boldsymbol{\theta}, \delta^*(T)).$$

It follows

$$R(\boldsymbol{\theta}, \delta) = \mathbb{E}_\theta[\mathbb{E}[L(\boldsymbol{\theta}, \delta(\mathbf{X}))|T]] \geq \mathbb{E}_\theta[L(\boldsymbol{\theta}, \mathbb{E}[\delta(\mathbf{X})|T])] = \mathbb{E}_\theta[L(\boldsymbol{\theta}, \delta^*(T))] = R(\boldsymbol{\theta}, \delta^*).$$

That is, δ^* is as good as δ . Therefore if δ is admissible, then δ^* is also admissible.

Minimax decision rules. Consider

$$\delta' \in \arg \min_{\delta \in \mathcal{D}} \left\{ \sup_{\boldsymbol{\theta} \in \Theta} R(\boldsymbol{\theta}, \delta) \right\}.$$

That is, decision rule δ' is minimax if

$$\sup_{\boldsymbol{\theta} \in \Theta} R(\boldsymbol{\theta}, \delta') = \inf_{\delta \in \mathcal{D}} \left\{ \sup_{\boldsymbol{\theta} \in \Theta} R(\boldsymbol{\theta}, \delta) \right\}.$$

Theorem 23.3 Suppose that δ is a unique minimax decision rule. Then δ is admissible.

Proof. Consider $\delta' \in \mathcal{D}$. Then since δ is minimax

$$\sup_{\theta \in \Theta} R(\theta, \delta') \geq \sup_{\theta \in \Theta} R(\theta, \delta).$$

Moreover since δ is unique we have, that if $\delta' \neq \delta$, then

$$\sup_{\theta \in \Theta} R(\theta, \delta') > \sup_{\theta \in \Theta} R(\theta, \delta),$$

i.e., δ' is not better than δ . It follows that δ is admissible. \square

How minimax decision rules are related to Bayes rules.

Proposition 23.1 *If*

$$\sup_{\theta \in \Theta} R(\theta, \delta^\pi) \leq B(\pi, \delta^\pi),$$

then δ^π is a minimax decision rule.

Proof. If δ^π is not minimax, then for some δ'

$$\sup_{\theta \in \Theta} R(\theta, \delta') < \sup_{\theta \in \Theta} R(\theta, \delta^\pi).$$

For any prior $\pi(\theta)$ we have

$$B(\pi, \delta) = \int R(\theta, \delta) \pi(\theta) d\theta \leq \sup_{\theta \in \Theta} R(\theta, \delta) \int \pi(\theta) d\theta = \sup_{\theta \in \Theta} R(\theta, \delta).$$

Hence then

$$B(\pi, \delta') \leq \sup_{\theta \in \Theta} R(\theta, \delta') < \sup_{\theta \in \Theta} R(\theta, \delta^\pi) \leq B(\pi, \delta^\pi),$$

which contradicts minimality of δ^π for $B(\pi, \cdot)$. \square

Saddle point

Consider problems

$$\max_{y \in Y} \min_{x \in X} g(x, y), \tag{23.3}$$

$$\min_{x \in X} \max_{y \in Y} g(x, y), \tag{23.4}$$

where X and Y are nonempty sets and $g : X \times Y \rightarrow \mathbb{R}$ is a real valued function. We have that for any $(x', y') \in X \times Y$,

$$\psi(y') = \min_{x \in X} g(x, y') \leq g(x', y') \leq \max_{y \in Y} g(x', y) = \varphi(x').$$

It follows that

$$\max_{y \in Y} \psi(y) \leq \min_{x \in X} \varphi(x).$$

Therefore we have that

$$\max_{y \in Y} \min_{x \in X} g(x, y) \leq \min_{x \in X} \max_{y \in Y} g(x, y), \tag{23.5}$$

i.e., optimal value of problem (23.3) is less then or equal to the optimal value of problem (23.4).

Now suppose that $\psi(\bar{y}) = \varphi(\bar{x})$ for some $(\bar{x}, \bar{y}) \in X \times Y$. By (23.5) this implies that optimal values of problems (23.3) and (23.4) are equal to each other and

$$\bar{y} \in \arg \max_{y \in Y} \psi(y) \quad \text{and} \quad \bar{x} \in \arg \min_{x \in X} \varphi(x).$$

That is

$$\max_{y \in Y} g(\bar{x}, y) = g(\bar{x}, \bar{y}) = \min_{x \in X} g(x, \bar{y}). \quad (23.6)$$

A point $(\bar{x}, \bar{y}) \in X \times Y$ satisfying the above condition (23.6) is called *saddle point*.

Let $(\bar{x}, \bar{y}) \in X \times Y$ be a saddle point. Then

$$\varphi(\bar{x}) = \max_{y \in Y} g(\bar{x}, y) = g(\bar{x}, \bar{y}) = \min_{x \in X} g(x, \bar{y}) = \psi(\bar{y}).$$

- It follows that if a saddle point (\bar{x}, \bar{y}) exists, then the optimal values of problems (23.3) and (23.4) are equal to each other, \bar{y} is an optimal solution of problem (23.3) and \bar{x} is an optimal solution of problem (23.4). Conversely if the optimal values of problems (23.3) and (23.4) are equal to each other, and \bar{y} is an optimal solution of problem (23.3) and \bar{x} is an optimal solution of problem (23.4), then (\bar{x}, \bar{y}) is a saddle point. \square

24 Spherical and elliptical distributions

An $m \times 1$ random vector \mathbf{X} is said to have *spherical* distribution if \mathbf{X} and $\mathbf{T}\mathbf{X}$ have the same distribution for any $m \times m$ orthogonal matrix \mathbf{T} .

Examples

- (i) Normal distribution $\mathbf{X} \sim \mathcal{N}_m(\mathbf{0}, \sigma^2 \mathbf{I}_m)$. The corresponding density function

$$f(\mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{m/2}} \exp\left(-\frac{1}{2}\sigma^{-2}\mathbf{x}'\mathbf{x}\right).$$

- (ii) ε -contaminated normal distribution, with pdf $(1 - \varepsilon)f_1(\mathbf{x}) + \varepsilon f_2(\mathbf{x})$, $\varepsilon \in [0, 1]$, where $f_i(\cdot)$ is pdf of $\mathcal{N}_m(\mathbf{0}, \sigma_i^2 \mathbf{I}_m)$, $i = 1, 2$.

- (iii) Multivariate t -distribution with n degrees of freedom. Its pdf is

$$f(\mathbf{x}) = \frac{\Gamma[\frac{1}{2}(n + m)]}{\Gamma(\frac{1}{2}n)(\pi n)^{m/2}} \frac{1}{(1 + n^{-1}\mathbf{x}'\mathbf{x})^{(n+m)/2}},$$

where $\Gamma(t) = \int_0^\infty x^{t-1} e^{-x} dx$. This is distribution of random vector $\mathbf{X} = Z^{-1/2} n^{1/2} \mathbf{Y}$, where $Z \sim \chi_n^2$ and $\mathbf{Y} \sim N_m(\mathbf{0}, \mathbf{I}_m)$, and Z and \mathbf{Y} are independent. This is the multivariate counterpart of t -distribution with n degrees of freedom.

Spherical distributions can be generated in the following way. Let X_1, \dots, X_m be random variables such that conditional on random variable $Z > 0$, $Z \sim G(\cdot)$, these variables are iid $N(0, Z)$. Then the pdf of random vector $\mathbf{X} = (X_1, \dots, X_m)'$ is

$$f(\mathbf{x}) = \int_0^\infty (2\pi z)^{-m/2} \exp\left(-\frac{1}{2}z^{-1}\mathbf{x}'\mathbf{x}\right) dG(z).$$

This is scale mixture of normal distributions. In particular, if Z can have two possible values σ_1^2 and σ_2^2 with respective probabilities $1 - \varepsilon$ and ε , then this is the ε -contaminated normal distribution. If $Z \sim n/\chi_n^2$, then \mathbf{X} has m -variate t -distribution with n degrees of freedom.

Recall that the characteristic function of a random vector \mathbf{X} is $\phi_X(\mathbf{t}) := \mathbb{E}[\exp(i\mathbf{t}'\mathbf{X})]$, where $i^2 = -1$ and $e^{i\theta} = \cos \theta + i \sin \theta$. If \mathbf{X} has spherical distribution, then \mathbf{X} and \mathbf{TX} have the same distribution for any orthogonal matrix \mathbf{T} and hence

$$\phi_X(\mathbf{t}) = \mathbb{E}[\exp(i\mathbf{t}'\mathbf{X})] = \mathbb{E}[\exp(i\mathbf{t}'\mathbf{TX})] = \mathbb{E}[\exp(i(\mathbf{T}'\mathbf{t})'\mathbf{X})] = \phi_X(\mathbf{T}'\mathbf{t}).$$

It follows that $\phi_X(\mathbf{t})$ is a function of $\mathbf{t}'\mathbf{t}$, i.e.,

$$\phi_X(\mathbf{t}) = \psi(\mathbf{t}'\mathbf{t}) \quad (24.1)$$

for some function $\psi(\cdot)$ of nonnegative real valued variable. Conversely suppose that the characteristic function of a random vector \mathbf{X} can be represented in the form (24.1). Then for any orthogonal matrix \mathbf{T} the characteristic function of \mathbf{X} is the same as the characteristic function of \mathbf{TX} and hence they have the same distribution. It follows that distribution of \mathbf{X} has spherical distribution iff the characteristic function of \mathbf{X} can be represented in the form (24.1).

It is said that an $m \times 1$ random vector \mathbf{X} has *elliptical* distribution with parameters $\boldsymbol{\mu} \in \mathbb{R}^m$ and symmetric positive definite $m \times m$ matrix $\mathbf{V} = [v_{ij}]_{i,j=1,\dots,m}$ if its pdf is

$$f(\mathbf{x}) = c_m |\mathbf{V}|^{-1/2} h((\mathbf{x} - \boldsymbol{\mu})'\mathbf{V}^{-1}(\mathbf{x} - \boldsymbol{\mu}))$$

for some function $h : \mathbb{R} \rightarrow \mathbb{R}_+$. The constant $c_m > 0$ is adjusted in such a way that $\int f(\mathbf{x})d\mathbf{x} = 1$. We use notation $\mathbf{X} \sim E_m(\boldsymbol{\mu}, \mathbf{V})$ for elliptical distributions. Note that $\mathbf{X} \sim E_m(\boldsymbol{\mu}, \mathbf{V})$ iff $\mathbf{Y} = \mathbf{V}^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$ has spherical distribution.

If $\mathbf{X} \sim E_m(\boldsymbol{\mu}, \mathbf{V})$, then $\mathbf{X} = \boldsymbol{\mu} + \mathbf{V}^{1/2}\mathbf{Y}$ where \mathbf{Y} has spherical distribution, and hence its characteristic function can be written as

$$\phi_X(\mathbf{t}) = \mathbb{E} \left[\exp(i\mathbf{t}'(\boldsymbol{\mu} + \mathbf{V}^{1/2}\mathbf{Y})) \right] = \exp(i\mathbf{t}'\boldsymbol{\mu})\mathbb{E}[\exp(i\mathbf{t}'\mathbf{V}^{1/2}\mathbf{Y})].$$

Since \mathbf{Y} has spherical distribution we have by (24.1) that

$$\mathbb{E}[\exp(i\mathbf{t}'\mathbf{V}^{1/2}\mathbf{Y})] = \psi((\mathbf{V}^{1/2}\mathbf{t})'(\mathbf{V}^{1/2}\mathbf{t})) = \psi(\mathbf{t}'\mathbf{V}\mathbf{t}).$$

That is, the characteristic function of $\mathbf{X} \sim E_m(\boldsymbol{\mu}, \mathbf{V})$ can be represented in the form

$$\phi_X(\mathbf{t}) = \exp(i\mathbf{t}'\boldsymbol{\mu})\psi(\mathbf{t}'\mathbf{V}\mathbf{t}) \quad (24.2)$$

for some function $\psi(\cdot)$. If $\mathbf{X} \sim N_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then by equation (2.2),

$$\phi_X(\mathbf{t}) = \exp(i\mathbf{t}'\boldsymbol{\mu} - \mathbf{t}'\boldsymbol{\Sigma}\mathbf{t}/2). \quad (24.3)$$

In that case $\psi(u) = e^{-u/2}$ for $u \geq 0$ with $\mathbf{V} = \boldsymbol{\Sigma}$.

Let $\mathbf{X} \sim E_m(\boldsymbol{\mu}, \mathbf{V})$ and \mathbf{A} be an $k \times m$ matrix of full row rank k . Then random vector $\mathbf{Y} = \mathbf{AX}$ has characteristic function

$$\phi_Y(\mathbf{t}) = \mathbb{E}[\exp(i\mathbf{t}'\mathbf{AX})] = \mathbb{E}[\exp(i(\mathbf{A}'\mathbf{t})'\mathbf{X})] = \phi_X(\mathbf{A}'\mathbf{t}) = \exp(i\mathbf{t}'\mathbf{A}\boldsymbol{\mu})\psi(\mathbf{t}'\mathbf{AV}\mathbf{A}'\mathbf{t}).$$

It follows that $\mathbf{Y} \sim E_k(\mathbf{A}\boldsymbol{\mu}, \mathbf{AV}\mathbf{A}')$. In particular let \mathbf{X} be partitioned $\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix}$, with the

corresponding partitioning of $\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}$ and $\mathbf{V} = \begin{bmatrix} \mathbf{V}_{11} & \mathbf{V}_{12} \\ \mathbf{V}_{21} & \mathbf{V}_{22} \end{bmatrix}$, where \mathbf{X}_1 is $m_1 \times 1$ and \mathbf{X}_2 is $m_2 \times 1$ subvectors of \mathbf{X} . Then $\mathbf{X}_1 \sim E_{m_1}(\boldsymbol{\mu}_1, \mathbf{V}_{11})$ with the characteristic function

$$\phi_{X_1}(\mathbf{t}_1) = \exp(i\mathbf{t}_1'\boldsymbol{\mu}_1)\psi(\mathbf{t}_1'\mathbf{V}_{11}\mathbf{t}_1), \quad (24.4)$$

and similarly for \mathbf{X}_2 .

Now suppose that components of random vector $\mathbf{X} = (X_1, \dots, X_m)$ have finite second order moments. Then

$$\partial \phi_{\mathbf{X}}(\mathbf{0}) / \partial \mathbf{t} = i \mathbb{E}[\mathbf{X}] \quad (24.5)$$

and

$$\partial^2 \phi_{\mathbf{X}}(\mathbf{0}) / \partial \mathbf{t} \partial \mathbf{t}' = -\mathbb{E}[\mathbf{X} \mathbf{X}'] = -\boldsymbol{\mu} \boldsymbol{\mu}' - \text{Cov}(\mathbf{X}). \quad (24.6)$$

It follows from (24.2) together with (24.5) and (24.6), that if $\mathbf{X} \sim E_m(\boldsymbol{\mu}, \mathbf{V})$, then $\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu}$ and $\text{Cov}(\mathbf{X}) = \alpha \mathbf{V}$, where $\alpha = -2\psi'(0)$. In particular this implies that

$$\text{Corr}(X_i, X_j) = \frac{v_{ij}}{\sqrt{v_{ii}v_{jj}}}, \quad i, j = 1, \dots, m. \quad (24.7)$$

By (24.7) we have that if $\mathbf{X} \sim E_m(\boldsymbol{\mu}, \mathbf{V})$ and $\mathbf{V} = \text{diag}(v_{11}, \dots, v_{mm})$ is diagonal, then X_1, \dots, X_m are uncorrelated.

Theorem 24.1 *Let $\mathbf{X} = (X_1, \dots, X_m)' \sim E_m(\boldsymbol{\mu}, \mathbf{V})$, $m \geq 2$, and $\mathbf{V} = \text{diag}(v_{11}, \dots, v_{mm})$. If X_1, \dots, X_m are all independent, then \mathbf{X} has multivariate normal distribution.*

Proof. By replacing \mathbf{X} with $\mathbf{X} - \boldsymbol{\mu}$, we can assume without loss of generality that $\boldsymbol{\mu} = \mathbf{0}$. Since matrix $\mathbf{V} = \text{diag}(v_{11}, \dots, v_{mm})$ is diagonal we have by (24.2) that the characteristic function of \mathbf{X} is

$$\phi_{\mathbf{X}}(\mathbf{t}) = \psi(\mathbf{t}' \mathbf{V} \mathbf{t}) = \psi\left(\sum_{i=1}^m t_i^2 v_{ii}\right).$$

Since X_1, \dots, X_m are independent, we have that

$$\phi_{\mathbf{X}}(\mathbf{t}) = \mathbb{E}[\exp(i\mathbf{t}' \mathbf{X})] = \mathbb{E}\left[\prod_{i=1}^m e^{it_i X_i}\right] = \prod_{i=1}^m \mathbb{E}\left[e^{it_i X_i}\right] = \prod_{i=1}^m \phi_i(t_i),$$

where $\phi_i(t_i)$ is the characteristic function of X_i , $i = 1, \dots, m$. By (24.4), $\phi_i(t_i) = \psi(t_i^2 v_{ii})$, and thus it follows that

$$\psi\left(\sum_{i=1}^m u_i^2\right) = \prod_{i=1}^m \psi(u_i^2), \quad (24.8)$$

where $u_i := t_i v_{ii}^{1/2}$.

In turn equation (24.8) implies that $\psi(u) = e^{-\kappa u^2}$ for some κ and $u > 0$. Indeed suppose that equation (24.8) holds. Then for any natural number p , $\psi(1) = \psi(1/p + \dots + 1/p) = \psi(1/p)^p$ and hence $\psi(1/p) = \psi(1)^{1/p}$. Furthermore for a rational positive number q/p we have $\psi(q/p) = \psi(1/p + \dots + 1/p) = \psi(1/p)^q$. It follows that $\psi(q/p) = \psi(1)^{q/p}$ for any positive rational number p/q . Moreover since function $\psi(\cdot)$ is continuous, it follows that for $u > 0$, $\psi(u) = \psi(1)^u = e^{-\kappa u^2}$ for $\kappa := -2 \log \psi(1)$. Since $\text{Cov}(\mathbf{X}) = \alpha \mathbf{V}$, where $\alpha = -2\psi'(0)$, we have then that $\text{Cov}(\mathbf{X}) = \kappa \mathbf{V}$, and hence $\kappa > 0$. It follows by (24.2) that the characteristic function of \mathbf{X} is $\phi_{\mathbf{X}}(\mathbf{t}) = \exp(-k\mathbf{t}' \mathbf{V} \mathbf{t}/2)$. That is, the characteristic function of \mathbf{X} coincides with the characteristic function of normal distribution with mean vector $\boldsymbol{\mu} = \mathbf{0}$ and covariance matrix $\boldsymbol{\Sigma} = k\mathbf{V}$. \square

24.1 Multivariate cumulants

Consider a random variable X . Let

$$\log \mathbb{E}[e^{tX}] = \sum_{n=1}^{\infty} \kappa_n \frac{t^n}{n!} \quad (24.9)$$

be Taylor expansion of its log-moments generating function (note that for $t = 0$ this function is 0). The coefficient κ_n is called n -th cumulant of X . Since $\mathbb{E}[e^{tX}]$ may not exist for $t \neq 0$, it is preferable to define cumulants in terms of the characteristic function as

$$\log \mathbb{E}[e^{itX}] = \sum_{n=1}^{\infty} \kappa_n \frac{(it)^n}{n!}, \quad (24.10)$$

where $\kappa_n = \left. \frac{\partial^n \log \mathbb{E}[e^{itX}]}{\partial t^n} \right|_{t=0}$.

Denote $\mu_k := \mathbb{E}[X^k]$ the k -th moment of X . Then

$$\begin{aligned} \kappa_1 &= \mu_1 = \mathbb{E}[X], \\ \kappa_2 &= \mu_2 - \mu_1^2 = \text{Var}(X), \\ \kappa_3 &= \mu_3 - 3\mu_1\mu_2 + 2\mu_1^3, \\ \kappa_4 &= \mu_4 - 4\mu_1\mu_3 - 3\mu_2^2 + 12\mu_2\mu_1^2 - 6\mu_1^4, \end{aligned}$$

provided these moments are finite. If X and Y are two independent random variables, then

$$\log \mathbb{E}[e^{it(X+Y)}] = \log \mathbb{E}[e^{itX}] + \log \mathbb{E}[e^{itY}],$$

and hence cumulants of $X + Y$ are equal to the sum of the respective cumulants of X and Y . In particular, if $Y = a$ where a is (deterministic) number, then the first cumulant of $X + a$ is $\kappa_1 + a$, and the cumulants of the higher order are the same as the cumulants of X .

Skewness of X is defined as

$$\gamma_1 := \frac{\kappa_3}{\kappa_2^{3/2}}, \quad (24.11)$$

kurtosis of X is defined as

$$\gamma_2 := \frac{\kappa_4}{\kappa_2^2}. \quad (24.12)$$

As it was pointed above, the skewness and kurtosis of X are the same as the respective skewness and kurtosis of $X + a$ for any number a . If distribution of X is symmetrical around its mean, then $\gamma_1 = 0$. If $X \sim N(0, \sigma^2)$, then $\mu_4 = 3\mu_2^2$ (see equation (2.4)). It follows that if $X \sim N(\mu, \sigma^2)$ then its kurtosis $\gamma_2 = 0$.

Consider now random vector $\mathbf{X} = (X_1, \dots, X_m)'$. Let $\phi_j(t_j)$ be the characteristic function of X_j . The cumulants of X_j are defined by

$$\log \phi_j(t_j) = \sum_{n=1}^{\infty} \kappa_n^j \frac{(it_j)^n}{n!}.$$

Mixed cumulants:

$$\log \phi_{j\ell}(t_j, t_\ell) = \sum_{n_1=1, n_2=1}^{\infty} \kappa_{n_1 n_2}^{j\ell} \frac{(it_j)^{n_1} (it_\ell)^{n_2}}{n_1! n_2!},$$

and so on.

Suppose that $\mathbf{X} \sim E_m(\boldsymbol{\mu}, \mathbf{V})$ has elliptical distribution. Then marginal distributions of X_j have zero skewness and the same kurtosis

$$\gamma_2^j = \frac{3[\psi''(0) - \psi'(0)^2]}{\psi'(0)^2}.$$

Denote $\kappa := \gamma_2^j/3$. Forth order cumulants of $\mathbf{X} \sim E_m(\boldsymbol{\mu}, \mathbf{V})$ are

$$\kappa_{1111}^{ijkl} = \kappa(\sigma_{ij}\sigma_{kl} + \sigma_{ik}\sigma_{jl} + \sigma_{il}\sigma_{jk}).$$

Let \mathbf{S} be the sample covariance matrix of sample of size N . By the CLT we have that $\mathbf{U}_N = N^{1/2}(\mathbf{S} - \boldsymbol{\Sigma})$ converges in distribution to normal with zero mean and covariances

$$\text{Cov}(u_{ij}, u_{kl}) = \kappa_{1111}^{ijkl} + \kappa_{11}^{ik}\kappa_{11}^{jl} + \kappa_{11}^{il}\kappa_{11}^{jk}.$$

If \mathbf{X} has normal distribution, then $\kappa = 0$ and

$$\text{Cov}(u_{ij}, u_{kl}) = \sigma_{ik}\sigma_{jl} + \sigma_{il}\sigma_{jk}.$$

Denote by $\boldsymbol{\Gamma}_N$ the corresponding $m^2 \times m^2$ covariance matrix (see section 15.2), where the subscript N emphasizes that this is under the assumption of normal distribution. For elliptical distribution, $N^{1/2}(\mathbf{s} - \boldsymbol{\sigma})$ converges in distribution to normal with zero mean and $m^2 \times m^2$ covariance matrix $\boldsymbol{\Gamma}$ with

$$\boldsymbol{\Gamma} = (1 + \kappa)\boldsymbol{\Gamma}_N + \kappa\boldsymbol{\sigma}\boldsymbol{\sigma}'.$$

25 Wishart distribution

Recall that

$$\mathbf{S} = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})'$$

is the sample covariance matrix of random sample $\mathbf{X}_1, \dots, \mathbf{X}_N$. Note that if $\mathbf{X}_1, \dots, \mathbf{X}_N$ is an iid sample from normal distribution $N_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then $\bar{\mathbf{X}}$ and \mathbf{S} are independent. Indeed

$$\text{Cov}(\bar{\mathbf{X}}, \mathbf{X}_i - \bar{\mathbf{X}}) = \text{Cov}(\bar{\mathbf{X}}, \mathbf{X}_i) - \text{Cov}(\bar{\mathbf{X}}).$$

Now $\text{Cov}(\bar{\mathbf{X}}, \mathbf{X}_i) = N^{-1}\boldsymbol{\Sigma}$ and $\text{Cov}(\bar{\mathbf{X}}) = N^{-2} \sum_{i=1}^N \text{Cov}(\mathbf{X}_i) = N^{-1}\boldsymbol{\Sigma}$. It follows that $\text{Cov}(\bar{\mathbf{X}}, \mathbf{X}_i - \bar{\mathbf{X}}) = \mathbf{0}$. That is, $\bar{\mathbf{X}}$ and $\mathbf{X}_i - \bar{\mathbf{X}}$ are uncorrelated, and because their joint distribution is normal, are independent. Since \mathbf{S} is a function of $\mathbf{X}_i - \bar{\mathbf{X}}$, $i = 1, \dots, N$, it follows that $\bar{\mathbf{X}}$ and \mathbf{S} are independent.

Let $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ be an iid sequence of random vectors having normal distribution $\mathcal{N}_m(\mathbf{0}, \boldsymbol{\Sigma})$. Consider random matrix

$$\mathbf{A} = \mathbf{Z}_1\mathbf{Z}_1' + \dots + \mathbf{Z}_n\mathbf{Z}_n'. \quad (25.1)$$

By definition \mathbf{A} has Wishart distribution, denoted $\mathbf{A} \sim W_m(n, \boldsymbol{\Sigma})$. In particular, for $m = 1$ and $Z_i \sim \mathcal{N}(0, \sigma^2)$, the corresponding A/σ^2 has chi-square distribution with n degrees of freedom.

Wishart distribution has the following properties.

- (i) If $\mathbf{A} \sim W_m(n, \boldsymbol{\Sigma})$ and $\alpha > 0$, then $\alpha\mathbf{A} \sim W_m(n, \alpha\boldsymbol{\Sigma})$. Indeed,

$$\alpha\mathbf{A} = (\alpha^{1/2}\mathbf{Z}_1)(\alpha^{1/2}\mathbf{Z}_1)' + \dots + (\alpha^{1/2}\mathbf{Z}_n)(\alpha^{1/2}\mathbf{Z}_n)',$$

and $\alpha^{1/2}\mathbf{Z}_i \sim \mathcal{N}(\mathbf{0}, \alpha\boldsymbol{\Sigma})$.

- (ii) If $\mathbf{A} \sim W_m(n, \boldsymbol{\Sigma})$ and \mathbf{B} is $m \times k$ deterministic matrix, then $\mathbf{B}'\mathbf{A}\mathbf{B} \sim W_k(n, \mathbf{B}'\boldsymbol{\Sigma}\mathbf{B})$. Indeed,

$$\mathbf{B}'\mathbf{A}\mathbf{B} = (\mathbf{B}'\mathbf{Z}_1)(\mathbf{B}'\mathbf{Z}_1)' + \dots + (\mathbf{B}'\mathbf{Z}_n)(\mathbf{B}'\mathbf{Z}_n)',$$

and $\mathbf{B}'\mathbf{Z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{B}'\boldsymbol{\Sigma}\mathbf{B})$.

(iii) Equation (25.1) can be written in the following form $\mathbf{A} = \mathbf{Z}'\mathbf{Z}$, where \mathbf{Z} is $n \times m$ matrix

$$\mathbf{Z} = \begin{bmatrix} \mathbf{z}'_1 \\ \cdot \\ \cdot \\ \mathbf{z}'_n \end{bmatrix} \text{ with } \mathbf{Z}' = [\mathbf{Z}_1, \dots, \mathbf{Z}_n]. \text{ Note that } \mathbb{E}[\mathbf{Z}] = \mathbf{0} \text{ and the covariance matrix of the corresponding } mn \times 1 \text{ vector } \text{vec}(\mathbf{Z}') \text{ is}^9$$

$$\text{Cov}(\text{vec}(\mathbf{Z}')) = \mathbf{I}_n \otimes \boldsymbol{\Sigma}. \quad (25.2)$$

Proposition 25.1 *Let $\mathbf{X}_1, \dots, \mathbf{X}_N \stackrel{iid}{\sim} \mathcal{N}_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and \mathbf{S} be the sample covariance matrix. Then $\mathbf{S} \sim W_m(n, n^{-1}\boldsymbol{\Sigma})$, where $n = N - 1$.*

Proof. Consider $N \times m$ matrix \mathbf{W} with rows $\mathbf{X}_i - \bar{\mathbf{X}}$, i.e., $\mathbf{W}' = [\mathbf{X}_1 - \bar{\mathbf{X}}, \dots, \mathbf{X}_N - \bar{\mathbf{X}}]$. Note that $\mathbf{S} = n^{-1}\mathbf{W}'\mathbf{W}$ and $\mathbf{W} = (\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N)\mathbf{X}$, where $\mathbf{X}' = [\mathbf{X}_1, \dots, \mathbf{X}_N]$. Matrix $\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N$ is a symmetric projection matrix of rank $N - 1$. Hence $\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N = \mathbf{H}\mathbf{H}'$, where \mathbf{H} is $N \times n$ matrix with $\mathbf{H}'\mathbf{H} = \mathbf{I}_n$ and $\mathbf{H}'\mathbf{1}_N = \mathbf{0}$ (spectral decomposition). Consider the following $n \times m$ matrix $\mathbf{Z} = \mathbf{H}'\mathbf{X}$. Then

$$\mathbf{Z}'\mathbf{Z} = \mathbf{X}'\mathbf{H}\mathbf{H}'\mathbf{X} = \mathbf{X}'(\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N)\mathbf{X} = \mathbf{W}'\mathbf{W},$$

and hence $\mathbf{S} = n^{-1}\mathbf{Z}'\mathbf{Z}$. Note that $\mathbb{E}[\mathbf{X}] = \mathbf{1}_N\boldsymbol{\mu}'$ and hence $\mathbb{E}[\mathbf{Z}] = \mathbf{H}'\mathbb{E}[\mathbf{X}] = \mathbf{H}'\mathbf{1}_N\boldsymbol{\mu}' = \mathbf{0}$. Now $\text{Cov}(\text{vec}(\mathbf{X}')) = \mathbf{I}_N \otimes \boldsymbol{\Sigma}$ and (see (15.22))

$$\text{vec}(\mathbf{X}'\mathbf{H}) = (\mathbf{H}' \otimes \mathbf{I}_m)\text{vec}(\mathbf{X}').$$

Thus using (15.21),

$$\text{Cov}(\text{vec}(\mathbf{Z}')) = \text{Cov}(\text{vec}(\mathbf{X}'\mathbf{H})) = (\mathbf{H}' \otimes \mathbf{I}_m)(\mathbf{I}_N \otimes \boldsymbol{\Sigma})(\mathbf{H} \otimes \mathbf{I}_m) = (\mathbf{H}'\mathbf{H}) \otimes \boldsymbol{\Sigma} = \mathbf{I}_n \otimes \boldsymbol{\Sigma}.$$

Hence $\mathbf{Z}'\mathbf{Z} \sim W_m(n, \boldsymbol{\Sigma})$, and $\mathbf{S} \sim W_m(n, n^{-1}\boldsymbol{\Sigma})$. \square

Theorem 25.1 *If $\mathbf{A} \sim W_m(n, \boldsymbol{\Sigma})$ and \mathbf{Y} is an $m \times 1$ random vector independent of \mathbf{A} and such that $\text{Prob}(\mathbf{Y} = \mathbf{0}) = 0$, then random variable $\frac{\mathbf{Y}'\mathbf{A}\mathbf{Y}}{\mathbf{Y}'\boldsymbol{\Sigma}\mathbf{Y}} \sim \chi_n^2$ and is independent of \mathbf{Y} .*

Proof. Conditional on \mathbf{Y} , we have that $\mathbf{Y}'\mathbf{A}\mathbf{Y} \sim W_1(n, \mathbf{Y}'\boldsymbol{\Sigma}\mathbf{Y})$ and hence

$$\frac{\mathbf{Y}'\mathbf{A}\mathbf{Y}}{\mathbf{Y}'\boldsymbol{\Sigma}\mathbf{Y}} \sim W_1(n, 1) = \chi_n^2.$$

That is, the conditional distribution of $\frac{\mathbf{Y}'\mathbf{A}\mathbf{Y}}{\mathbf{Y}'\boldsymbol{\Sigma}\mathbf{Y}}$ does not depend on \mathbf{Y} . It follows that $\frac{\mathbf{Y}'\mathbf{A}\mathbf{Y}}{\mathbf{Y}'\boldsymbol{\Sigma}\mathbf{Y}}$ is independent of \mathbf{Y} and its (unconditional) distribution is χ_n^2 . \square

Together with Proposition 25.1 this implies the following.

Proposition 25.2 *Let $\mathbf{X}_1, \dots, \mathbf{X}_N \stackrel{iid}{\sim} \mathcal{N}_m(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and \mathbf{S} be the corresponding sample covariance matrix. Then $\frac{n\bar{\mathbf{X}}'\mathbf{S}\bar{\mathbf{X}}}{\bar{\mathbf{X}}'\boldsymbol{\Sigma}\bar{\mathbf{X}}} \sim \chi_n^2$ and is independent of $\bar{\mathbf{X}}$ (recall that $n = N - 1$).*

Proof. Since \mathbf{S} and $\bar{\mathbf{X}}$ are independent and $n\mathbf{S} \sim W_m(n, \boldsymbol{\Sigma})$, the result follows from Theorem 25.1. \square

⁹Recall definitions of Kronecker product of matrices and vec operator discussed in section 15.2.

Theorem 25.2 Let $\mathbf{A} \sim W_m(n, \boldsymbol{\Sigma})$ be partitioned $\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}$, where \mathbf{A}_{11} is of order $k \times k$ and \mathbf{A}_{22} is of order $(m - k) \times (m - k)$, and matrix $\boldsymbol{\Sigma}$ is partitioned accordingly $\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$. Consider $\mathbf{A}_{11.2} = \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}$ and $\boldsymbol{\Sigma}_{11.2} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}$. Then

$$\mathbf{A}_{11.2} \sim W_k(n - m + k, \boldsymbol{\Sigma}_{11.2}), \quad (25.3)$$

and $\mathbf{A}_{11.2}$ is independent of \mathbf{A}_{22} .

Proof. Since $\mathbf{A} \sim W_m(n, \boldsymbol{\Sigma})$ it can be written in the form (25.1), or equivalently as $\mathbf{A} = \mathbf{Z}'\mathbf{Z}$, where $\mathbf{Z} = \begin{bmatrix} \mathbf{z}'_1 \\ \cdot \\ \cdot \\ \cdot \\ \mathbf{z}'_n \end{bmatrix}$ is the respective $n \times m$ matrix. Let us partition $\mathbf{Z} = [\tilde{\mathbf{Z}}_1, \tilde{\mathbf{Z}}_2]$, where

$\tilde{\mathbf{Z}}_1$ is of order $n \times k$ and $\tilde{\mathbf{Z}}_2$ is of order $n \times (m - k)$. Note that $\tilde{\mathbf{Z}}_1 = \begin{bmatrix} \mathbf{z}'_{11} \\ \cdot \\ \cdot \\ \cdot \\ \mathbf{z}'_{1n} \end{bmatrix}$ and $\tilde{\mathbf{Z}}_2 = \begin{bmatrix} \mathbf{z}'_{21} \\ \cdot \\ \cdot \\ \cdot \\ \mathbf{z}'_{2n} \end{bmatrix}$,

where $\mathbf{Z}_i = \begin{bmatrix} \mathbf{z}_{1i} \\ \mathbf{z}_{2i} \end{bmatrix}$ are respective partitions of vectors \mathbf{Z}_i , $i = 1, \dots, n$. Recall that conditional on $\mathbf{Z}_{2i} = \mathbf{z}_2$,

$$\mathbf{Z}_{1i} \sim \mathcal{N}(\boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\mathbf{z}_2, \boldsymbol{\Sigma}_{11.2}) \quad (25.4)$$

(see equation (2.3)). Note that matrix $\mathbf{I}_n - \tilde{\mathbf{Z}}_2(\tilde{\mathbf{Z}}_2'\tilde{\mathbf{Z}}_2)^{-1}\tilde{\mathbf{Z}}_2'$ is idempotent (projection) of rank $n - (m - k) = n - m + k$, and

$$[\mathbf{I}_n - \tilde{\mathbf{Z}}_2(\tilde{\mathbf{Z}}_2'\tilde{\mathbf{Z}}_2)^{-1}\tilde{\mathbf{Z}}_2']\tilde{\mathbf{Z}}_2 = \mathbf{0}. \quad (25.5)$$

Because of (25.4) and (25.5) we have that conditional on $\tilde{\mathbf{Z}}_2$,

$$\tilde{\mathbf{Z}}_1'[\mathbf{I}_n - \tilde{\mathbf{Z}}_2(\tilde{\mathbf{Z}}_2'\tilde{\mathbf{Z}}_2)^{-1}\tilde{\mathbf{Z}}_2']\tilde{\mathbf{Z}}_1 \sim W_k(n - m + k, \boldsymbol{\Sigma}_{11.2}).$$

Moreover

$$\tilde{\mathbf{Z}}_1'[\mathbf{I}_n - \tilde{\mathbf{Z}}_2(\tilde{\mathbf{Z}}_2'\tilde{\mathbf{Z}}_2)^{-1}\tilde{\mathbf{Z}}_2']\tilde{\mathbf{Z}}_1 = \underbrace{\tilde{\mathbf{Z}}_1'\tilde{\mathbf{Z}}_1}_{\mathbf{A}_{11}} - \underbrace{\tilde{\mathbf{Z}}_1'\tilde{\mathbf{Z}}_2}_{\mathbf{A}_{12}} \underbrace{(\tilde{\mathbf{Z}}_2'\tilde{\mathbf{Z}}_2)^{-1}}_{\mathbf{A}_{22}^{-1}} \underbrace{\tilde{\mathbf{Z}}_2'\tilde{\mathbf{Z}}_1}_{\mathbf{A}_{21}} = \mathbf{A}_{11.2}.$$

It follows that the (unconditional) distribution of $\mathbf{A}_{11.2}$ is $W_k(n - m + k, \boldsymbol{\Sigma}_{11.2})$, and that $\mathbf{A}_{11.2}$ is independent of $\tilde{\mathbf{Z}}_2$ and hence of \mathbf{A}_{22} . \square

Theorem 25.3 Let $\mathbf{A} \sim W_m(n, \boldsymbol{\Sigma})$ and \mathbf{B} be (deterministic) $m \times k$ matrix of rank k . Then $(\mathbf{B}'\mathbf{A}^{-1}\mathbf{B})^{-1} \sim W_k(n - m + k, (\mathbf{B}'\boldsymbol{\Sigma}^{-1}\mathbf{B})^{-1})$.

Proof. Note that the assertion is invariant under linear transformations. That is, if \mathbf{C} is an $m \times m$ nonsingular matrix, then by replacing \mathbf{B} with $\tilde{\mathbf{B}} = \mathbf{C}\mathbf{B}$ and \mathbf{A} with $\tilde{\mathbf{A}} = \mathbf{C}\mathbf{A}\mathbf{C}'$ we have $\tilde{\mathbf{B}}'\tilde{\mathbf{A}}^{-1}\tilde{\mathbf{B}} = \mathbf{B}'\mathbf{A}^{-1}\mathbf{B}$. Moreover $\tilde{\mathbf{A}} \sim W_m(n, \tilde{\boldsymbol{\Sigma}})$, where $\tilde{\boldsymbol{\Sigma}} = \mathbf{C}\boldsymbol{\Sigma}\mathbf{C}'$, and $\tilde{\mathbf{B}}'\tilde{\boldsymbol{\Sigma}}^{-1}\tilde{\mathbf{B}} = \mathbf{B}'\boldsymbol{\Sigma}^{-1}\mathbf{B}$.

Therefore by applying an appropriate linear transformation, we can assume that $\mathbf{B} = \begin{bmatrix} \mathbf{I}_k \\ \mathbf{0} \end{bmatrix}$.

Then $\mathbf{B}'\mathbf{A}^{-1}\mathbf{B} = \mathbf{A}^{11}$, where $\mathbf{A}^{-1} = \begin{bmatrix} \mathbf{A}^{11} & \mathbf{A}^{12} \\ \mathbf{A}^{21} & \mathbf{A}^{22} \end{bmatrix}$. Now (see (2.7))

$$\mathbf{A}^{11} = (\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})^{-1},$$

and hence $(\mathbf{B}'\mathbf{A}^{-1}\mathbf{B})^{-1} = \mathbf{A}_{11.2}$. By Theorem 25.2 we have that $\mathbf{A}_{11.2} \sim W_k(n - m + k, \Sigma_{11.2})$. It remains to note that here $\Sigma_{11.2} = (\mathbf{B}'\Sigma^{-1}\mathbf{B})^{-1}$. \square

Proposition 25.3 *If $\mathbf{A} \sim W_m(n, \Sigma)$ and \mathbf{Y} is an $m \times 1$ random vector independent of \mathbf{A} and such that $\text{Prob}(\mathbf{Y} = \mathbf{0}) = 0$, then*

$$\frac{\mathbf{Y}'\Sigma^{-1}\mathbf{Y}}{\mathbf{Y}'\mathbf{A}^{-1}\mathbf{Y}} \sim \chi_{n-m+1}^2. \quad (25.6)$$

Proof. By Theorem 25.3 we have that conditional on \mathbf{Y} , $(\mathbf{Y}'\mathbf{A}^{-1}\mathbf{Y})^{-1} \sim W_1(n - m + 1, (\mathbf{Y}'\Sigma^{-1}\mathbf{Y})^{-1})$. This implies (25.6). \square

25.1 Hotelling's T^2 statistic

Hotelling's T^2 statistic is an extension of t distribution to a multivariate setting. Let $\mathbf{X}_1, \dots, \mathbf{X}_N$ be an iid sample from normal distribution $\mathcal{N}_m(\boldsymbol{\mu}, \Sigma)$, and \mathbf{S} be the sample covariance matrix. Recall that $\bar{\mathbf{X}}$ and \mathbf{S} are independent.

Suppose that we want to test $H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$ against $H_1 : \boldsymbol{\mu} \neq \boldsymbol{\mu}_0$, where $\boldsymbol{\mu}_0$ is a given $m \times 1$ vector. Hotelling's T^2 statistic for testing H_0 is

$$T^2 = N(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)' \mathbf{S}^{-1}(\bar{\mathbf{X}} - \boldsymbol{\mu}_0). \quad (25.7)$$

For $m = 1$ this statistic can be written as $\frac{(\bar{X} - \mu_0)^2}{S^2/N}$, where $S^2 = (N - 1)^{-1} \sum_{i=1}^N (X_i - \bar{X})^2$ is the sample variance. So in that case $T^2 = t^2$, where $t = \frac{\bar{X} - \mu_0}{S/\sqrt{N}}$ is the usual t statistic.

We proceed now to statistical inference of Hotelling's statistic. For $n = N - 1$ we can write

$$\frac{T^2}{n} = \frac{N(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)' \Sigma^{-1}(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)}{\left(\frac{n(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)' \Sigma^{-1}(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)}{(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)' \mathbf{S}^{-1}(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)} \right)}$$

Under H_0 we have that $N^{1/2}(\bar{\mathbf{X}} - \boldsymbol{\mu}_0) \sim \mathcal{N}(\mathbf{0}, \Sigma)$, and hence $N(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)' \Sigma^{-1}(\bar{\mathbf{X}} - \boldsymbol{\mu}_0) \sim \chi_m^2$. Also by Proposition 25.1 we have that $n\mathbf{S} \sim W_m(n, \Sigma)$ and hence by Proposition 25.3,

$$\frac{n(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)' \Sigma^{-1}(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)}{(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)' \mathbf{S}^{-1}(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)} \sim \chi_{n-m+1}^2.$$

We obtain the following result.

Theorem 25.4 *Let $\mathbf{X}_1, \dots, \mathbf{X}_N \stackrel{iid}{\sim} N_m(\boldsymbol{\mu}, \Sigma)$. Then under $H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$,*

$$\frac{(N - m)T^2}{m(N - 1)} \sim F_{m, N-m}. \quad (25.8)$$

Note that as $N \rightarrow \infty$, the coefficient $\frac{N-m}{m(N-1)}$ in (25.8) tends to $1/m$. Therefore for large N the distribution of T^2 becomes like χ_m^2 . This should be not surprising since by the LLN, \mathbf{S} converges w.p.1 to Σ , and $N(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)' \Sigma^{-1}(\bar{\mathbf{X}} - \boldsymbol{\mu}_0)$ has χ_m^2 distribution when $\boldsymbol{\mu} = \boldsymbol{\mu}_0$ (Theorem 3.1).

Suppose now that we want to test linear model $H_0 : \mathbf{A}\boldsymbol{\mu} = \mathbf{c}$, where \mathbf{A} is a $k \times m$ matrix of rank k and \mathbf{c} is $k \times 1$ vector. The corresponding Hotelling's T^2 statistic is

$$T^2 = N \min_{\mathbf{A}\boldsymbol{\mu}=\mathbf{c}} (\bar{\mathbf{X}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}). \quad (25.9)$$

It is possible to write this in the form

$$T^2 = N(\mathbf{A}\bar{\mathbf{X}} - \mathbf{c})'(\mathbf{A}\mathbf{S}\mathbf{A}')^{-1}(\mathbf{A}\bar{\mathbf{X}} - \mathbf{c}). \quad (25.10)$$

Indeed, suppose for the sake of simplicity that $\mathbf{c} = \mathbf{0}$. Consider $\tilde{\mathbf{X}} = \mathbf{S}^{-1/2}\bar{\mathbf{X}}$ and $\tilde{\mathbf{A}} = \mathbf{A}\mathbf{S}^{1/2}$. Then making change of variables $\boldsymbol{\tau} = \mathbf{S}^{-1/2}\boldsymbol{\mu}$ we have

$$\min_{\mathbf{A}\boldsymbol{\mu}=\mathbf{0}} (\bar{\mathbf{X}} - \boldsymbol{\mu})' \mathbf{S}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}) = \min_{\tilde{\mathbf{A}}\boldsymbol{\tau}=\mathbf{0}} (\tilde{\mathbf{X}} - \boldsymbol{\tau})' (\tilde{\mathbf{X}} - \boldsymbol{\tau}). \quad (25.11)$$

The right hand side of (25.11) is the squared distance from $\tilde{\mathbf{X}}$ to the space orthogonal to the one generated by matrix $\tilde{\mathbf{A}}$. Hence

$$\min_{\tilde{\mathbf{A}}\boldsymbol{\tau}=\mathbf{0}} (\tilde{\mathbf{X}} - \boldsymbol{\tau})' (\tilde{\mathbf{X}} - \boldsymbol{\tau}) = \tilde{\mathbf{X}}' \tilde{\mathbf{A}}' (\tilde{\mathbf{A}}\tilde{\mathbf{A}}')^{-1} \tilde{\mathbf{A}}\tilde{\mathbf{X}} = (\mathbf{A}\bar{\mathbf{X}})' (\mathbf{A}\mathbf{S}\mathbf{A}')^{-1} (\mathbf{A}\bar{\mathbf{X}}).$$

Under H_0 , $N^{1/2}(\mathbf{A}\bar{\mathbf{X}} - \mathbf{c}) \sim \mathcal{N}_k(\mathbf{0}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}')$ and

$$\frac{(N-k)T^2}{k(N-1)} \sim F_{k, N-k}. \quad (25.12)$$

Indeed, consider $\mathbf{Y}_i = \mathbf{A}\mathbf{X}_i$, $i = 1, \dots, N$. We have that $\mathbf{Y}_i \sim N_k(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}')$. Also $\bar{\mathbf{Y}} = \mathbf{A}\bar{\mathbf{X}}$ and the corresponding sample covariance matrix is $\mathbf{A}\mathbf{S}\mathbf{A}'$. Hotelling's T^2 statistic for testing $H_0 : \mathbf{A}\boldsymbol{\mu} = \mathbf{c}$ is given by the left hand side of (25.12).

26 Spatial statistics

Consider a (real valued) function $Z(\mathbf{x})$ of $x \in \mathbb{R}^d$. Given values (observations, measurements) of $Z(\cdot)$ at some points, we would like to evaluate (to estimate) value of $Z(\mathbf{x})$ at a given point $x = x^*$. As a modeling approach we view $Z(\mathbf{x})$ as a random process. It is said that $Z(\mathbf{x})$ is stationary if for any points $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^d$ and $\mathbf{h} \in \mathbb{R}^d$, random vector $(Z(\mathbf{x}_1), \dots, Z(\mathbf{x}_m))$ has the same distribution as $(Z(\mathbf{x}_1 + \mathbf{h}), \dots, Z(\mathbf{x}_m + \mathbf{h}))$. This definition of stationarity is too general for practical use. It is said that $Z(\mathbf{x})$ is second order (or weakly) stationary if its mean $\mathbb{E}[Z(\mathbf{x})]$ is constant (independent of \mathbf{x}), and its covariance function $c(\mathbf{x}, \mathbf{y}) = \text{Cov}(Z(\mathbf{x}), Z(\mathbf{y}))$ has the property that for any $\mathbf{x}, \mathbf{y}, \mathbf{h} \in \mathbb{R}^d$ it follows that $c(\mathbf{x} + \mathbf{h}, \mathbf{y} + \mathbf{h}) = c(\mathbf{x}, \mathbf{y})$. Of course any stationary process is second order stationary provided it has finite second order moments. By taking $\mathbf{h} = -\mathbf{y}$ we have then that $c(\mathbf{x}, \mathbf{y}) = c(\mathbf{x} - \mathbf{y}, \mathbf{0})$. That is, for the second order stationary process the covariance function depends on the difference $\mathbf{x} - \mathbf{y}$. So we use notation $c(\mathbf{x} - \mathbf{y}) = \text{Cov}(Z(\mathbf{x}), Z(\mathbf{y}))$ for the (auto)covariance function.

The autocovariance function $c(\cdot)$ has the following properties. It is symmetric, i.e., $c(\mathbf{h}) = c(-\mathbf{h})$, this follows from that $\text{Cov}(Z(\mathbf{x}), Z(\mathbf{y})) = \text{Cov}(Z(\mathbf{y}), Z(\mathbf{x}))$. Since $c(\mathbf{0}) = \text{Cov}(Z(\mathbf{x}), Z(\mathbf{x})) = \text{Var}(Z(\mathbf{x}))$, it follows that $c(\mathbf{0}) > 0$. We have that

$$|\text{Cov}(Z(\mathbf{x}), Z(\mathbf{y}))| \leq \sqrt{\text{Var}(Z(\mathbf{x}))} \sqrt{\text{Var}(Z(\mathbf{y}))}$$

and hence $|c(\mathbf{h})| \leq c(\mathbf{0})$ for all $\mathbf{h} \in \mathbb{R}^d$. The function $c(\cdot)$ should be positive definite. That is for any $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^d$ the covariance matrix of $(Z(\mathbf{x}_1), \dots, Z(\mathbf{x}_m))$ should be positive semidefinite, i.e., the $m \times m$ matrix with entries $a_{ij} = c(\mathbf{x}_i - \mathbf{x}_j)$, $i, j = 1, \dots, m$, should be positive semidefinite.

The semivariogram of (stationary) process $Z(\mathbf{x})$ is defined as

$$\gamma(\mathbf{h}) := \frac{1}{2}\mathbb{E}[|Z(\mathbf{x} + \mathbf{h}) - Z(\mathbf{x})|^2].$$

Note that we can assume that $\mathbb{E}[Z(\mathbf{h})] = 0$ and hence $c(\mathbf{0}) = \text{Var}(Z(\mathbf{h})) = \mathbb{E}[Z(\mathbf{h})^2]$, and thus

$$\gamma(\mathbf{h}) = \frac{1}{2}\mathbb{E}[Z(\mathbf{x} + \mathbf{h})^2 + Z(\mathbf{x})^2 - 2Z(\mathbf{x} + \mathbf{h})Z(\mathbf{x})] = c(\mathbf{0}) - c(\mathbf{h}).$$

Consider $m \times m$ matrix $\mathbf{\Gamma}$ with entries $\Gamma_{ij} = \gamma(\mathbf{x}_i - \mathbf{x}_j)$, $i, j = 1, \dots, m$. Note that $\Gamma_{ij} = c_0 - c_{ij}$, where $c_0 = c(\mathbf{0})$ and $c_{ij} = c(\mathbf{x}_i - \mathbf{x}_j)$. In matrix form this can be written as $\mathbf{\Gamma} = c_0\mathbf{1}_m\mathbf{1}'_m - \mathbf{C}$, where \mathbf{C} is $m \times m$ matrix with entries c_{ij} .

Given observations $Z(\mathbf{x}_1), \dots, Z(\mathbf{x}_N)$ consider the linear predictor

$$\hat{Z}(\mathbf{x}) = \sum_{i=1}^N w_i Z(\mathbf{x}_i).$$

We have that

$$\mathbb{E}[\hat{Z}(\mathbf{x})] = \sum_{i=1}^N w_i \mathbb{E}[Z(\mathbf{x}_i)] = \mu \sum_{i=1}^N w_i,$$

where μ is the mean of the process. Therefore $\hat{Z}(\mathbf{x})$ is unbiased iff $\sum_{i=1}^N w_i = 1$. It is said that $\hat{Z}(\mathbf{x})$ is the Best Linear Unbiased Predictor (BLUP) if the weights w_i are chosen to minimize variance of the error $\hat{Z}(\mathbf{x}) - Z(\mathbf{x})$. Now (since $\sum_{i=1}^N w_i = 1$)

$$\text{Var}(\hat{Z}(\mathbf{x}) - Z(\mathbf{x})) = \text{Var} \left[\sum_{i=1}^N w_i (Z(\mathbf{x}_i) - Z(\mathbf{x})) \right],$$

and

$$\text{Cov}(Z(\mathbf{x}_i) - Z(\mathbf{x}), Z(\mathbf{x}_j) - Z(\mathbf{x})) = \text{Cov}(Z(\mathbf{x}_i), Z(\mathbf{x}_j)) - \text{Cov}(Z(\mathbf{x}), Z(\mathbf{x}_i)) - \text{Cov}(Z(\mathbf{x}), Z(\mathbf{x}_j)) + c(\mathbf{0}).$$

Moreover

$$\text{Cov}(Z(\mathbf{x}_i), Z(\mathbf{x}_j)) = c(\mathbf{0}) - \gamma(\mathbf{x}_i - \mathbf{x}_j) = c_0 - \Gamma_{ij}.$$

In matrix form we can write this as

$$\text{Var}(\hat{Z}(\mathbf{x}) - Z(\mathbf{x})) = -\mathbf{w}'\mathbf{\Gamma}\mathbf{w} + 2\mathbf{g}'\mathbf{w},$$

where $\Gamma_{ij} = \gamma(\mathbf{x}_i - \mathbf{x}_j)$ and $g_i = \gamma(\mathbf{x} - \mathbf{x}_i)$. The BLUP is solution of the problem

$$\min_{\mathbf{w}} -\mathbf{w}'\mathbf{\Gamma}\mathbf{w} + 2\mathbf{g}'\mathbf{w} \quad \text{subject to} \quad \sum_{i=1}^N w_i = 1.$$

By using method of Lagrange multipliers this can be written as the following system of $N + 1$ linear equations

$$\begin{bmatrix} \gamma(\mathbf{x}_1 - \mathbf{x}_1) & \cdots & \gamma(\mathbf{x}_1 - \mathbf{x}_N) & 1 \\ \cdots & \cdots & \cdots & \cdots \\ \gamma(\mathbf{x}_N - \mathbf{x}_1) & \cdots & \gamma(\mathbf{x}_N - \mathbf{x}_N) & 1 \\ 1 & \cdots & 1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ \cdots \\ w_N \\ \lambda \end{bmatrix} = \begin{bmatrix} \gamma(\mathbf{x} - \mathbf{x}_1) \\ \cdots \\ \gamma(\mathbf{x} - \mathbf{x}_N) \\ 1 \end{bmatrix}$$

with $N + 1$ unknowns w_1, \dots, w_N, λ .

It is said that the model is isotropic if $\gamma(\mathbf{h})$ is a function of $\|\mathbf{h}\|$. In that case the semivariogram $\gamma(h)$ becomes a function of one dimensional variable $h = \|\mathbf{h}\|$. The following are some popular parametric models of semivariograms.

Linear $\gamma(0) = 0$ and $\gamma(h) = c_0 + bh$ for $h > 0$, where $c_0 \geq 0$ and $b > 0$ are parameters. This model is valid for any dimension d . Note that here $\lim_{h \downarrow 0} \gamma(h) = c_0$ with c_0 could be strictly positive. Value $\lim_{h \downarrow 0} \gamma(h)$ is called the nugget effect.

Exponential model $\gamma(0) = 0$ and $\gamma(h) = c_0 + c_\ell(1 - e^{-h/a_\ell})$ for $h > 0$, where $c_0 \geq 0$, $c_\ell > 0$ and $a_\ell > 0$. This model is valid for any dimension d .

Note that both models have nugget c_0 , and in the linear model the semivariogram is unbounded, while in the exponential model the semivariogram is bounded by $c_0 + c_\ell$.

Positive-definite functions. Recall that for complex number $c = a + bi$ its conjugate $\bar{c} = a - bi$, where $i^2 = -1$. A function $\phi : \mathbb{R}^n \rightarrow \mathbb{C}$ is positive-definite if $\phi(-\mathbf{x}) = \overline{\phi(\mathbf{x})}$ and for any $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$ and $c_1, \dots, c_m \in \mathbb{C}$ it follows that

$$\sum_{k,\ell=1}^m c_k \bar{c}_\ell \phi(\mathbf{x}_k - \mathbf{x}_\ell) \geq 0.$$

This means that the corresponding $m \times m$ matrix $\mathbf{\Gamma}$ with components $\gamma_{k\ell} = \phi(\mathbf{x}_k - \mathbf{x}_\ell)$ is Hermitian¹⁰. If $\phi(\mathbf{x})$ is real valued, then the corresponding matrix $\mathbf{\Gamma}$ is positive semidefinite. For $m = 1$ and $c_1 = 1$ it follows that $\phi(\mathbf{x}_1 - \mathbf{x}_1) \geq 0$, i.e., $\phi(\mathbf{0}) \geq 0$. Also $|\phi(\mathbf{x})| \leq \phi(\mathbf{0})$.

Recall that $e^{i\theta} = \cos \theta + i \sin \theta$. Consider Fourier transform of finite positive Borel measure μ on \mathbb{R}^n

$$\hat{\mu}(\mathbf{z}) = \int_{\mathbb{R}^n} e^{-i\mathbf{z}'\mathbf{x}} d\mu(\mathbf{x}), \quad \mathbf{z} \in \mathbb{R}^n.$$

If $d\mu(\mathbf{x}) = f(\mathbf{x})d\mathbf{x}$, then

$$\hat{\mu}(\mathbf{z}) = \int_{\mathbb{R}^n} e^{-i\mathbf{z}'\mathbf{x}} f(\mathbf{x})d\mathbf{x},$$

is the Fourier transform of function f . Note that measure μ is positive if $f(\mathbf{x}) \geq 0$ for all $\mathbf{x} \in \mathbb{R}^n$.

For any $\mathbf{z}_1, \dots, \mathbf{z}_m \in \mathbb{R}^n$ and $c_1, \dots, c_m \in \mathbb{C}$ we have

$$\begin{aligned} \sum_{k,\ell=1}^m c_k \bar{c}_\ell \hat{\mu}(\mathbf{z}_k - \mathbf{z}_\ell) &= \sum_{k,\ell=1}^m c_k \bar{c}_\ell \int_{\mathbb{R}^n} e^{-i(\mathbf{z}_k - \mathbf{z}_\ell)'\mathbf{x}} d\mu(\mathbf{x}) = \\ \int_{\mathbb{R}^n} \sum_{k,\ell=1}^m c_k \bar{c}_\ell e^{-i(\mathbf{z}_k - \mathbf{z}_\ell)'\mathbf{x}} d\mu(\mathbf{x}) &= \int_{\mathbb{R}^n} \left(\sum_{k=1}^m c_k e^{-i\mathbf{z}'_k \mathbf{x}} \right) \overline{\left(\sum_{\ell=1}^m c_\ell e^{-i\mathbf{z}'_\ell \mathbf{x}} \right)} d\mu(\mathbf{x}) = \\ \int_{\mathbb{R}^n} \left| \sum_{k=1}^m c_k e^{-i\mathbf{z}'_k \mathbf{x}} \right|^2 d\mu(\mathbf{x}) &\geq 0. \end{aligned}$$

That is, Fourier transform of a finite positive Borel measure is a positive definite function. The converse of that is also true (its proof is not trivial).

Theorem 26.1 (Bochner) *If $\phi : \mathbb{R}^n \rightarrow \mathbb{C}$ is positive definite, continuous, and satisfies $\phi(\mathbf{0}) = 1$, then there is Borel probability measure μ on \mathbb{R}^n such that ϕ is Fourier transform of μ .*

¹⁰A matrix $\mathbf{A} = [a_{k\ell}]$ is said to be Hermitian if $a_{k\ell} = \bar{a}_{\ell k}$ and $\sum_{k,\ell=1}^m a_{k\ell} x_k \bar{x}_\ell \geq 0$ for any $x_1, \dots, x_m \in \mathbb{C}$.