

TMA4267 - Linear statistical models

Trond Skaret Johansen

Spring semester 2024

Contents

1	Multivariate Distribution and its Generalisations	4
1.1	Matrix algebra	4
1.2	Random vectors and distributions	5
1.3	Multivariate expecations and moments	6
1.4	Transformations	7
1.4.1	Mahalanobis transformation	7
1.4.2	Principal component analysis (PCA)	7
1.4.3	General transformations	9
1.5	Characteristic function	9
2	The Multivariate Normal Distribution	11
2.1	From univariate to multivariate	11
2.1.1	Univariate case	11
2.1.2	Multivariate case	11
2.2	Estimation of the multivariate normal distribution	13
2.2.1	Univariate case	13
2.2.2	Multivariate case	14
2.2.3	Quadratic forms	15
2.2.4	Idempotent matrices	15
3	Multiple Linear Regression	17
3.1	Model and assumptions	17
3.2	Parameter estimation	17
3.2.1	Properties of the the estimators, fitted values and residuals	18
3.2.2	Inference about β_j	19
3.3	Some notes on independence	20
3.4	Analysis of variance (ANOVA)	20
3.4.1	Fictional model	21
3.4.2	Further expressions for the sums of squares	22
3.5	F-test	22
3.6	General F-test	22
3.7	Transformations of data	23
3.7.1	Box-Cox transformation	24

3.7.2	Variance stabilising transformation	24
4	Model Analysis, Selection and Multiple Hypothesis Testing	25
4.1	Model analysis	25
4.2	Model selection	25
4.2.1	Which model to choose	25
4.3	Multiple hypothesis testing	26
4.3.1	p -value is a random variable	26
4.3.2	Testing m hypothesis	27
5	ANOVA and Design of Experiment	29
5.1	Analysis of variance (ANOVA)	29
5.1.1	Two-way ANOVA	30
5.2	Two level factorial design	30
5.2.1	Modelling interactions	31
5.2.2	Inference about effect	31
5.3	Fractional factorial design (2^{k-r} -design)	32
5.3.1	Blocking	33

Introduction

This is a brief summary of the course TMA4267 about linear statistical models. It includes the main content from the lecture held by **TODO**: ... recorded in, where some examples etc... are excluded.

The purpose of the notes is to give a good overview of the syllabus. I intend to add summaries of the lectures as I review them. I hope to include insights from projects / exercises where it is appropriate.

Topics

The first chapter begins by introducing multivariate distributions and how to compute expectations. We then move on to multivariate moments and transformations. Principal component analysis (PCA) is described, before we explain the need for characteristic functions and not just moment generating functions when working with multivariate distributions.

In the second chapter we introduce the multivariate normal distribution. We deal with estimation in the multivariate normal distribution and give the theory of quadratic forms and idempotent matrices.

The third chapter tackles multiple linear regression. We introduce the model and its assumptions and estimate the parameters. Properties of the estimators, fitted values and residuals are established, and then put to use in performing inference about the coefficients. We perform t-tests, do ANOVA, compute the coefficient of determination and perform general F-tests. Finally, we look at some way of transforming the data.

In the fourth chapter, we analyse the model and the selection of models. We perform multiple hypothesis testing and present examples. Two methods for controlling the FWER are investigated.

In the fifth and final chapter, we do more ANOVA and investigate design of experiment. Our focus is two-level factorial design. We discuss generators of design, resolution of design and the method of blocking.

Course progress

- | | | |
|-----------------|-----------------|---------------|
| • First reading | ✓ Lecture 11-12 | ✓ June 2019 |
| ✓ Lecture 1-25 | ✓ Lecture 13-14 | ✓ May 2018 |
| • Gjennomgang | □ Lecture 15-16 | ✓ May 2017 |
| ✓ Lecture 1-2 | ✓ Lecture 17-18 | ✓ June 2016 |
| ✓ Lecture 3-4 | ✓ Lecture 19-20 | □ May 2015 |
| ✓ Lecture 5-6 | ✓ Lecture 21-22 | □ May 2014 |
| ✓ Lecture 7-8 | • Exams | □ August 2014 |
| ✓ Lecture 9-10 | ✓ May 2023 | |

1 Multivariate Distribution and its Generalisations

1.1 Matrix algebra

The main tool of the course is matrix algebra. We therefore repeat the most important concepts. The matrix \mathbf{A} is said to be *symmetric* if $\mathbf{A}^\top = \mathbf{A}$. It is *orthogonal* if $\mathbf{A}\mathbf{A}^\top = \mathbf{A}^\top\mathbf{A} = \mathbf{I}$, i.e. if the columns are orthogonal. The elements of the pair $(\lambda, \boldsymbol{\gamma})$ are called *eigenvalue* and *eigenvector* respectively if they satisfy $\mathbf{A}\boldsymbol{\gamma} = \lambda\boldsymbol{\gamma}$ and $\boldsymbol{\gamma} \neq \mathbf{0}$. λ can be found as the solution of $\det(\mathbf{A} - \lambda\mathbf{I})$. Recall that the *trace*, $\text{tr}(\mathbf{A})$, of a matrix is the sum of the diagonal. We also have the formulas:

$$\det \mathbf{A} = \prod_{i=1}^p \lambda_i \quad \text{tr}(\mathbf{A}) = \sum_{i=1}^p \lambda_i.$$

For a symmetric matrix we may find the *Jordan decomposition*. Let $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_p\}$ and $\Gamma = (\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_p)$. Then we have:

$$\mathbf{A} = \Gamma \Lambda \Gamma^\top.$$

For a symmetric matrix \mathbf{A} and a vector \mathbf{x} we define the *quadratic form*:

$$Q(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} = \sum_{i=1}^p \sum_{j=1}^p x_i A_{ij} x_j.$$

Theorem 1. *Transforming $\mathbf{y} = \Gamma^\top \mathbf{x}$ we obtain*

$$Q(\mathbf{x}) = \sum_{i=1}^p \lambda_i y_i^2.$$

A matrix is said to be *positive definite* if $Q(\mathbf{x}) > 0$ for all $\mathbf{x} \neq 0$ and positive semi definite if $Q(\mathbf{x}) \geq 0$ for all $\mathbf{x} \neq 0$. We write $A > 0$ and $A \geq 0$ respectively.

Theorem 2. *The symmetric matrix A is positive definite iff $\lambda_i > 0$ for all i .*

Proof. Using the transform of the previous theorem we find

$$\lambda_1 y_1^2 + \dots + \lambda_p y_p^2 > 0 \quad \forall \mathbf{y} \in \mathbb{R}^p \Leftrightarrow \lambda_i > 0 \quad \forall i.$$

□

From this we obtain two more useful results.

1. If $A > 0$ the inverse exists and the determinant is > 0
2. If $A > 0$ there exists a unique positive definite square root with decomposition:

$$A^{1/2} = \Gamma \Lambda^{1/2} \Gamma^\top.$$

1.2 Random vectors and distributions

A *random vector* \mathbf{X} is a vector where each component is a *random variable*. Similarly, we define a *random matrix* as a matrix with random variables as component. As in the univariate case, we define the *cumulative distribution function* (CDF) by:

$$F(\mathbf{x}) = \mathbb{P}[\mathbf{X} \leq \mathbf{x}] = \mathbb{P}[X_1 \leq x_1, \dots, X_p \leq x_p].$$

A random vector is said to be *absolutely continuous* if there exists a *probability density function* (PDF) f such that:

$$F(\mathbf{x}) = \int_{-\infty}^{x_p} \dots \int_{-\infty}^{x_1} f(u_1, \dots, u_p) du_1 \dots du_p.$$

Then, we may compute the probability of the event $\mathbf{X} \in D$ by:

$$\mathbb{P}[\mathbf{X} \in D] = \int_D f(\mathbf{x}) d\mathbf{x} \quad \forall D \subseteq \mathbb{R}^p.$$

The random vector is said to be *discrete* if it is concentrated on a countable (finite or infinite) set of points. Then integral becomes a sum.

In the absolutely continuous case, we may write the density as:

$$f(\mathbf{x}) = f(x_1, \dots, x_p) = \frac{\partial^p F(x_1, \dots, x_p)}{\partial x_1 \dots \partial x_p}.$$

Let $\mathbf{X}_A, \mathbf{X}_B$ be two random vectors st. $\mathbf{X} = (\mathbf{X}_A, \mathbf{X}_B)^\top$ has cdf F . Then we may find the *marginal distribution*:

$$F_A(x_1, \dots, x_k) = F(x_1, \dots, x_k, \infty, \dots, \infty).$$

In absolutely continuous case we find the marginal density:

$$f_A(x_1, \dots, x_k) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(x_1, \dots, x_p) du_p \dots du_{k+1}$$

The *conditional distribution* of \mathbf{X}_B given $\mathbf{X}_A = \mathbf{x}_A$ is:

$$f_{\mathbf{X}_B|\mathbf{X}_A=\mathbf{x}_A}(\mathbf{x}_B) = \frac{f(x_1, \dots, x_p)}{f_A(x_1, \dots, x_k)}$$

We say that the random vectors $\mathbf{X}_A, \mathbf{X}_B$ are *independent* if

$$F(x_1, \dots, x_p) = F_A(x_1, \dots, x_k) F_B(x_{k+1}, \dots, x_p) \quad \forall x_1, \dots, x_p.$$

In the continuous case we have independence iff $f = f_A \cdot f_B$. In this case $f_{\mathbf{x}_B|\mathbf{X}_A=\mathbf{x}_A} = f_B(\mathbf{x}_B)$. Similar definition for independence when \mathbf{X} has N components and not just 2.

1.3 Multivariate expectations and moments

We define the *expectation* of the random vector \mathbf{X} as

$$\mathbb{E}[\mathbf{X}] = (\mathbb{E}[X_1], \dots, \mathbb{E}[X_p])^\top.$$

Here, each $\mathbb{E}[X_i] = \int_{\mathbb{R}^p} x_i f(\mathbf{x}) d\mathbf{x}$. It is easy to show:

1. For constants a, b we have $\mathbb{E}[a\mathbf{X} + b\mathbf{Y}] = a\mathbb{E}[\mathbf{X}] + b\mathbb{E}[\mathbf{Y}]$.
2. For (shape compatible) matrices \mathbf{A}, \mathbf{B} we have $\mathbb{E}[\mathbf{A}\mathbf{X}\mathbf{B}] = \mathbf{A}\mathbb{E}[\mathbf{X}]\mathbf{B}$.
3. Let \mathbf{X}, \mathbf{Y} be *independent* random matrices whose product is defined. Then $\mathbb{E}[\mathbf{X}\mathbf{Y}] = \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{Y}]$.

Let $\mathbf{X} = (X_1, \dots, X_p)^\top$ and $\mathbb{E}[\mathbf{X}] =: \boldsymbol{\mu} = (\mu_1, \dots, \mu_p)$. We then define the *covariance matrix* of \mathbf{X} as:

$$\text{Var}[\mathbf{X}] = \text{Cov}[\mathbf{X}] = \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{X_1 X_1} & \cdots & \sigma_{X_1 X_p} \\ \vdots & \ddots & \vdots \\ \sigma_{X_p X_1} & \cdots & \sigma_{X_p X_p} \end{pmatrix} = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^\top].$$

This matrix is *symmetric*. Note also that $\Sigma_{ij} = \text{Cov}[X_i, X_j]$. We can also show:

$$\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{X}\mathbf{X}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top.$$

The correlation matrix (with ones on the diagonal as $\rho_{X_i X_i} = 1$) is given by

$$\boldsymbol{\rho} = \begin{pmatrix} \rho_{X_1 X_1} & \cdots & \rho_{X_1 X_p} \\ \vdots & \ddots & \vdots \\ \rho_{X_p X_1} & \cdots & \rho_{X_p X_p} \end{pmatrix}, \quad \rho_{X_i X_j} = \frac{\sigma_{X_i X_j}}{\sqrt{\sigma_{X_i}} \sqrt{\sigma_{X_j}}}.$$

For two random vectors \mathbf{X}, \mathbf{Y} we define their covariance matrix by

$$\boldsymbol{\Sigma}_{\mathbf{XY}} = \text{Cov}[\mathbf{X}, \mathbf{Y}] = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu}_{\mathbf{X}})(\mathbf{Y} - \boldsymbol{\mu}_{\mathbf{Y}})^\top] = (\text{Cov}[X_i, X_j])_{\substack{i=1, \dots, p \\ j=1, \dots, q}}$$

Proposition 1. *The covariance matrix $\boldsymbol{\Sigma}$ is positive semi-definite.*

Proof. Using the formula for the variance of a linear combination we obtain:

$$\begin{aligned} \mathbf{y}^\top \boldsymbol{\Sigma} \mathbf{y} &= (y_1 \quad \cdots \quad y_n) \begin{pmatrix} \Sigma_{11} & \Sigma_{12} & \cdots & \Sigma_{1n} \\ \Sigma_{21} & \Sigma_{22} & \cdots & \Sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \Sigma_{n1} & \Sigma_{n2} & \cdots & \Sigma_{nn} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \sum_{i=1}^n \sum_{j=1}^n y_i \Sigma_{ij} y_j \\ &= \sum_{i=1}^n y_i^2 \text{Var}(X_i) + 2 \sum_{i < j} y_i y_j \text{Cov}(X_i, X_j) = \text{Var} \left(\sum_{i=1}^n y_i X_i \right) \geq 0. \end{aligned}$$

Which completes the proof. □

Remark. We usually require that the covariance matrix is [positive definite](#), since if it is only positive semi-definite there are nontrivial linear combinations with 0 variance. Indeed, if $\Sigma = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ we obtain using the formula for variance of linear combinations that:

$$\text{Var}[X_1 - X_2] = 1^2 \text{Var}[X_1] + (-1)^2 \text{Var}[X_2] + 2(1)(-1) \text{Cov}[X_1, X_2] = 0.$$

This never happens for positive definite matrices since we for any nontrivial linear combination have:

$$\text{Cov}[\mathbf{c}^\top \mathbf{X}] = \mathbf{c}^\top \Sigma \mathbf{c} > 0.$$

We have many more properties of covariance matrices:

1. $\Sigma_{\mathbf{X}\mathbf{Y}} = \Sigma_{\mathbf{Y}\mathbf{X}}^\top$.
2. If $\mathbf{X} \sim (\boldsymbol{\mu}_\mathbf{X}, \Sigma_{\mathbf{X}\mathbf{X}})$, $\mathbf{Y} \sim (\boldsymbol{\mu}_\mathbf{Y}, \Sigma_{\mathbf{Y}\mathbf{Y}})$ then $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})^\top$ has

$$\Sigma_{\mathbf{Z}\mathbf{Z}} = \begin{pmatrix} \Sigma_{\mathbf{X}\mathbf{X}} & \Sigma_{\mathbf{X}\mathbf{Y}} \\ \Sigma_{\mathbf{Y}\mathbf{X}} & \Sigma_{\mathbf{Y}\mathbf{Y}} \end{pmatrix}.$$

3. Independence of \mathbf{X}, \mathbf{Y} implies $\text{Cov}[\mathbf{X}, \mathbf{Y}] = \mathbf{0}$ (NB: the converse not true).
4. $\text{Var}[\mathbf{A}\mathbf{X} + \mathbf{b}] = \mathbf{A} \text{Var}[\mathbf{X}] \mathbf{A}^\top$.
5. $\text{Cov}[\mathbf{X} + \mathbf{Y}, \mathbf{Z}] = \text{Cov}[\mathbf{X}, \mathbf{Z}] + \text{Cov}[\mathbf{Y}, \mathbf{Z}]$.
6. $\text{Var}[\mathbf{X} + \mathbf{Y}] = \text{Var}[\mathbf{X}] + \text{Cov}[\mathbf{X}, \mathbf{Y}] + \text{Cov}[\mathbf{Y}, \mathbf{X}] + \text{Var}[\mathbf{Y}]$.
7. $\text{Cov}[\mathbf{A}\mathbf{X}, \mathbf{B}\mathbf{Y}] = \mathbf{A} \text{Cov}[\mathbf{X}, \mathbf{Y}] \mathbf{B}^\top$.

1.4 Transformations

1.4.1 Mahalanobis transformation

Our first transformation is the [Mahalanobis transformation](#). We recall that we can get 0 mean and unit variance in the univariate case by the transformation $Y = \frac{X - \mu}{\sigma}$. In the multivariate case, suppose $\mathbf{X} = (X_1, \dots, X_p)^\top \sim (\boldsymbol{\mu}, \Sigma)$ with Σ non-singular. Then using the unique positive definite square root $\Sigma^{1/2}$ of Σ , we have the transformation:

$$\boxed{\mathbf{Y} = \Sigma^{-1/2}(\mathbf{X} - \boldsymbol{\mu}) \sim (0, \mathbf{I})} \tag{1}$$

Proof. Compute the expected value and variance. □

1.4.2 Principal component analysis (PCA)

In the following we suppose that we observe realisations of some random vector $\mathbf{X} = (X_1, \dots, X_p)^\top \sim (\boldsymbol{\mu}, \Sigma)$. The goal of [principal component analysis](#) is to reduce dimensionality by removing some components and keeping components with large variance and hence more information.

The idea is to first perform the transform $\mathbf{X} \mapsto \mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b} = (Y_1, \dots, Y_p)^\top$ in such a way that:

1. $\mathbb{E}[\mathbf{Y}] = \mathbf{0}$,

2. $\text{Cov}[Y_i, Y_j] = 0$ for $i \neq j$, i.e. Σ_Y is diagonal,
3. $\text{Var}[Y_1] \geq \text{Var}[Y_2] \geq \dots \geq \text{Var}[Y_p]$.

Using again the Jordan decomposition we may find such a transform. Arrange $\Sigma = \Gamma \Lambda \Gamma^\top$ such that the eigenvalues on the diagonal of $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_p\}$ are ordered decreasingly. A transformation satisfying the requirements is then $Y = \Gamma^{1/2}(\mathbf{X} - \mu)$. Indeed we may compute $\text{Var}[Y] = \Lambda$. We call Y_1, \dots, Y_p the *principal components*.

For the interpretation, recall that for the trace of a matrix we have: $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA}) = \text{tr}(\mathbf{CAB})$. Using this we compute:

$$\text{“total variance”} = \sum_{i=1}^p \text{Var}[X_i] = \text{tr}(\Sigma) = \text{tr}(\Gamma \Lambda \Gamma^\top) = \text{tr}(\Lambda \Gamma^\top \Gamma) = \text{tr}(\Lambda) = \sum_{i=1}^n \lambda_i.$$

Since we care about variability in our data, we can use this interpretation to extract the features explaining as much variance as possible. For example we can keep the m first principal components. Say we compute:

$$\frac{\sum_{i=1}^m \lambda_i}{\sum_{i=1}^p \lambda_i} = 0.8.$$

Then we may say that “80% of variability is explained by the components”, and we have successfully reduced dimensionality.

We can also do *empirical PCA* (also called *sample PCA*). Suppose that we have n observations of independent identically distributed random vectors $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^\top \sim N(\mu, \Sigma)$. Gather them in a $(n \times p)$ “data matrix”:

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1^\top \\ \mathbf{X}_2^\top \\ \vdots \\ \mathbf{X}_n^\top \end{pmatrix} = \begin{pmatrix} X_{11} & \dots & X_{1p} \\ X_{21} & \dots & X_{2p} \\ \vdots & \ddots & \vdots \\ X_{n1} & \dots & X_{np} \end{pmatrix}.$$

Estimate the unknown μ, Σ by $\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$ and $\mathbf{S} = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{X}_j - \bar{\mathbf{X}})(\mathbf{X}_j - \bar{\mathbf{X}})^\top$. Then we may transform the matrix \mathbf{X} into $\mathbf{Y} = (\mathbf{X} - \mathbf{1}_n \bar{\mathbf{X}}) \mathbf{G}$ where $\mathbf{S} = \mathbf{G} \mathbf{L} \mathbf{G}$ is the Jordan decomposition of \mathbf{S} with eigenvalues ordered decreasingly. Finally remove lowest components.

Remark. Intuitively, PCA can be thought of as finding a p -dimensional ellipsoid that fits the data. The axis of the ellipsoid represent the principal components.

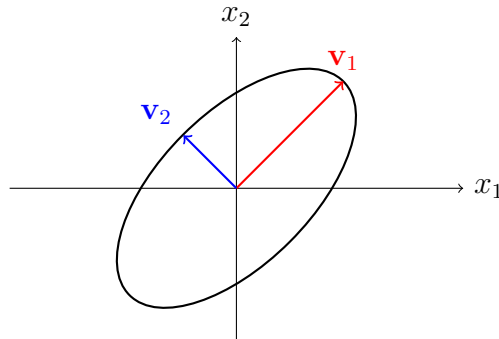


Figure 1: Principal Component Analysis (PCA) Ellipsoid. The red arrow indicates the first principal component (eigenvectors) \mathbf{v}_1 and the blue arrow indicates the second principal component \mathbf{v}_2 .

1.4.3 General transformations

Suppose we transform the random vector $\mathbf{X} = (X_1, \dots, X_p)^\top$ into $\mathbf{Y} = g(\mathbf{X})$. If the function g is one-to-one and has differentiable inverse u , then

$$f_{\mathbf{Y}}(\mathbf{y}) = |\det \mathbf{J}| f_{\mathbf{X}}(u(\mathbf{y}))$$

where

$$\mathbf{J} = \left(\frac{\partial u_i(y)}{\partial y_j} \right)_{i,j=1,\dots,p}$$

is the *Jacobian matrix*. An important special case is that of linear transformations $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b}$. If \mathbf{A} is non-singular the inverse transform is $\mathbf{X} = \mathbf{A}^{-1}(\mathbf{Y} - \mathbf{b})$ and $\mathbf{J} = \mathbf{A}^{-1}$, so we obtain:

$$f_{\mathbf{Y}}(\mathbf{y}) = |\det \mathbf{A}^{-1}| f_{\mathbf{X}}(\mathbf{A}^{-1}(\mathbf{y} - \mathbf{b})).$$

1.5 Characteristic function

In the univariate case we often studied the *moment generating function*:

$$M_X(t) = \mathbb{E} [e^{tX}].$$

Important properties include:

1. $M_X = M_Y \Rightarrow X \stackrel{d}{=} Y$,
2. $\mathbb{E} [X^k] = M_X^{(k)}(0)$,
3. Independence of X, Y implies $M_{X+Y} = M_X M_Y$.

However, it does not exist for for instance the student-t distribution. It only exists when all moments exist. We shall now define the *characteristic function*, which **always** exists:

$$\varphi_X(t) = \mathbb{E} [e^{itX}].$$

It has similar properties:

1. $\varphi_X = \varphi_Y \Rightarrow X \stackrel{d}{=} Y$,
2. $\mathbb{E} [X^k] = i^{-k} \varphi^{(k)}(0)$,
3. Independence of X, Y implies $\varphi_{X+Y} = \varphi_X \varphi_Y$.

In the p-variate case, the functions are functions of $\mathbf{t} = (t_1, \dots, t_p)^\top$. Let as usual $\mathbf{X} = (X_1, \dots, X_p)^\top$ be our random vector. We define:

$$M_{\mathbf{X}}(\mathbf{t}) = \mathbb{E} [e^{\mathbf{t}^\top \mathbf{X}}],$$

$$\varphi_{\mathbf{X}}(\mathbf{t}) = \mathbb{E} [e^{i\mathbf{t}^\top \mathbf{X}}].$$

We list some important properties.

1. If $\varphi_{\mathbf{X}}(t)$ is absolutely integrable, the (Fourier) inversion formula holds:

$$f_{\mathbf{X}} = \frac{1}{(2\pi)^p} \int_{\mathbb{R}} e^{-i\mathbf{t}^\top \mathbf{x}} \varphi_{\mathbf{X}}(t) dt.$$

2. Denote $t_{(1)} = (t_1, \dots, 0)^\top, \dots, t_{(p)} = (0, \dots, t_p)^\top$. Then

$$\varphi_{\mathbf{X}_k}(t_k) = \mathbb{E} [e^{it_1 X_1}] = \mathbb{E} [e^{it_{(1)}^\top \mathbf{X}}] = \varphi_{\mathbf{X}}(t_{(k)}).$$

3. Let $\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}$ and $\mathbf{t} = \begin{pmatrix} \mathbf{t}_1 \\ \mathbf{t}_2 \end{pmatrix}$ be vectors with appropriate dimensions. Then:

$$\boxed{\mathbf{X}_1, \mathbf{X}_2 \text{ are independent} \Leftrightarrow \varphi_{\mathbf{X}}(\mathbf{t}) = \varphi_{\mathbf{X}_1}(\mathbf{t}_1)\varphi_{\mathbf{X}_2}(\mathbf{t}_2)} \quad (2)$$

4. Let \mathbf{X}, \mathbf{Y} be independent p-variate random vectors. Then:

$$\varphi_{\mathbf{X}+\mathbf{Y}}(\mathbf{t}) = \varphi_{\mathbf{X}}(\mathbf{t})\varphi_{\mathbf{Y}}(\mathbf{t}).$$

We end the section with a theorem linking the distributions of 1D random variables to the distribution of the p-variate random vector \mathbf{X} :

Theorem 3. [Cramer-Wold](#) theorem states that the distribution of $\mathbf{X} \in \mathbb{R}^p$ is completely determined by the set of all 1D distributions of $\mathbf{t}^\top \mathbf{X}, \mathbf{t} \in \mathbb{R}^p$.

Proof. Suppose that for all $\mathbf{t} \in \mathbb{R}^p$ we have $\mathbf{t}^\top \mathbf{X} \stackrel{d}{=} \mathbf{t}^\top \mathbf{Y}$. Then:

$$\mathbf{t}^\top \mathbf{X} \stackrel{d}{=} \mathbf{t}^\top \mathbf{Y} \Rightarrow \varphi_{\mathbf{t}^\top \mathbf{X}}(u) = \varphi_{\mathbf{t}^\top \mathbf{Y}}(u).$$

Taking $u = 1$ gives $\mathbb{E} [e^{it^\top \mathbf{X}}] = \mathbb{E} [e^{it^\top \mathbf{Y}}]$ for all \mathbf{t} , so $\varphi_{\mathbf{X}} = \varphi_{\mathbf{Y}}$ and hence $\mathbf{X} \stackrel{d}{=} \mathbf{Y}$. □

2 The Multivariate Normal Distribution

2.1 From univariate to multivariate

2.1.1 Univariate case

We begin by recalling important properties and results of the univariate normal distribution. Let $-\infty < \mu < \infty, \sigma \geq 0$. We say $X \sim N(\mu, \sigma^2)$ if it has density:

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

The case $\sigma = 0$ is considered normal degenerate and we have $X \stackrel{\text{a.s.}}{=} \mu$. One can show $\mathbb{E}[X] = \mu$, $\text{Var}[X] = \sigma^2$ and that the moment generating and characteristic functions are:

$$\begin{aligned} M_X(t) &= e^{\mu t + \frac{\sigma^2 t^2}{2}}, \\ \varphi_X(t) &= e^{i\mu t - \frac{\sigma^2 t^2}{2}}. \end{aligned}$$

We say that $Z \sim N(0, 1)$ is standard normal. We can transform as follows:

$$\begin{aligned} X \sim N(\mu, \sigma^2) &\Rightarrow \frac{X - \mu}{\sigma} \sim N(0, 1), \\ Z \sim N(0, 1) &\Rightarrow \sigma Z + \mu \sim N(\mu, \sigma^2). \end{aligned}$$

A final fundamental property is that linear combinations of independent normal random variables are also normal.

2.1.2 Multivariate case

We choose to present one of several equivalent definitions. Here we first define a special case and then generalise it.

Definition 1. A p -variate random vector $\mathbf{Z} = (Z_1, \dots, Z_p)^\top$ is a *standard normal vector* if its components are independent and each $Z_i \sim N(0, 1)$. We immediately obtain $f_{\mathbf{Z}}(\mathbf{z}) = \prod f_{Z_i}(z_i)$. We use the notation $\mathbf{Z} \sim N(0, \mathbf{I})$.

Definition 2. We say that $\mathbf{X} = (X_1, \dots, X_p)^\top$ is *multivariate normal* if there exists a $(p \times q)$ matrix \mathbf{A} and $\boldsymbol{\mu} \in \mathbb{R}^p$ such that

$$\mathbf{X} = \mathbf{A}\mathbf{Z} + \boldsymbol{\mu}, \quad \mathbf{Z} \sim (0, \mathbf{I}_q)$$

Theorem 4. For any vector $\mathbf{b} \in \mathbb{R}^p$ we have that $\mathbf{b}^\top \mathbf{X}$ is a univariate normal random variable.

Proof. It is a linear combination of independent normal random variables:

$$\mathbf{b}^\top \mathbf{X} = \mathbf{b}^\top (\mathbf{A}\mathbf{Z} + \boldsymbol{\mu}) = \mathbf{b}^\top \mathbf{A}\mathbf{Z} + \mathbf{b}^\top \boldsymbol{\mu} \sim N.$$

□

From the definition we find $\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu}$, $\text{Var}[\mathbf{X}] = \mathbf{A}\mathbf{A}^\top =: \boldsymbol{\Sigma}$. We write $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. It turns out that the statement from the theorem would give an equivalent definition!

Theorem 5. Let $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. If $\boldsymbol{\Sigma} > 0$, then the probability density function (PDF) of \mathbf{X} is:

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{(2\pi)^{k/2} \det \boldsymbol{\Sigma}^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})} \quad (3)$$

Proof. Recall that $\mathbf{X} = \boldsymbol{\Sigma}^{1/2} \mathbf{Z} + \boldsymbol{\mu}$. Using this, the result follows from the transformation formula and using the pdf of the standard normal \mathbf{Z} . \square

Theorem 6.

$$\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \Rightarrow \mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top) \quad (4)$$

Proof. Compute expected value and covariance. Normality follows from definition as we may write $\mathbf{Y} = \mathbf{A}(\mathbf{B}\mathbf{Z} + \mathbf{c}) + \mathbf{b} = \tilde{\mathbf{B}}\mathbf{Z} + \tilde{\mathbf{b}}$. \square

The following corollary tells us that we may *standardize* in the multivariate case as well:

Corollary 1.

$$\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \Rightarrow \mathbf{Z} = \boldsymbol{\Sigma}^{-1/2}(\mathbf{X} - \boldsymbol{\mu}) \sim N(\mathbf{0}, \mathbf{I}) \quad (5)$$

Corollary 2. Any subvector \mathbf{X}^* of $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is also normal.

Proof. Take \mathbf{A} to be a “component removing” matrix and use Theorem 6. \square

Theorem 7. The characteristic function of $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is:

$$\varphi_{\mathbf{X}}(\mathbf{t}) = e^{i\mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2}\mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t}} \quad (6)$$

Proof. For $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I}_p)$, we may use independence to obtain:

$$\varphi_{\mathbf{Z}}(\mathbf{t}) = \prod_{i=1}^p \varphi_{Z_i}(t_i) = e^{-\frac{1}{2}\mathbf{t}^\top \mathbf{t}}.$$

Then write $\mathbf{X} = \boldsymbol{\Sigma}^{1/2} \mathbf{Z} + \boldsymbol{\mu}$ and rewrite:

$$\varphi_{\mathbf{X}}(\mathbf{t}) = \mathbb{E} \left[e^{i\mathbf{t}^\top \mathbf{X}} \right] = e^{i\mathbf{t}^\top \boldsymbol{\mu}} \mathbb{E} \left[e^{i\mathbf{t}^\top \boldsymbol{\Sigma}^{1/2} \mathbf{Z}} \right] = e^{i\mathbf{t}^\top} \varphi_{\mathbf{Z}}(\boldsymbol{\Sigma}^{1/2} \mathbf{t}) = e^{i\mathbf{t}^\top} e^{-\frac{1}{2}\mathbf{t}^\top \boldsymbol{\Sigma} \mathbf{t}}.$$

Where we have used symmetry of $\boldsymbol{\Sigma}^{1/2}$ a few times. \square

The following is among the most important results in the course:

Theorem 8. Suppose $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Then

$$\mathbf{A}\mathbf{X}, \mathbf{B}\mathbf{X} \text{ are independent} \Leftrightarrow \text{Cov}[\mathbf{A}\mathbf{X}, \mathbf{B}\mathbf{X}] = \mathbf{A}\boldsymbol{\Sigma}\mathbf{B}^\top = \mathbf{0} \quad (7)$$

Proof. The rightward implication is true for any random vectors. To prove the leftward for normal \mathbf{X} we show that we may factor the characteristic function

$$\varphi \begin{pmatrix} \mathbf{A}\mathbf{X} \\ \mathbf{B}\mathbf{X} \end{pmatrix}(\mathbf{t}) = \varphi_{\mathbf{A}\mathbf{X}}(\mathbf{t}_1)\varphi_{\mathbf{B}\mathbf{X}}(\mathbf{t}_2).$$

□

Corollary 3. Let $\mathbf{X} = (X_1, \dots, X_p)^\top \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Then X_i, X_j independent iff $\text{Cov}[X_i, X_j] = 0$.

The next theorem tells us that we may “remove the dependent part” of a component from another:

Theorem 9. Let $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ be p -variate and $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2)^\top$ with $\mathbf{X}_i \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$. Define $\mathbf{X}'_2 = \mathbf{X}_2 - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\mathbf{X}_1$. Then:

1. $\mathbf{X}'_2 \sim N(\boldsymbol{\mu}_2 - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12})$.
2. $\mathbf{X}_1, \mathbf{X}'_2$ are independent.

Theorem 10. With notation as in the previous theorem, the conditional distribution of \mathbf{X}_2 given $\mathbf{X}_1 = \mathbf{x}_1$ is:

$$\mathbf{X}_2|\mathbf{X}_1=\mathbf{x}_1 \sim N(\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\mathbf{x}_1 - \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12}).$$

Theorem 11. If $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and $\boldsymbol{\Sigma}$ is non-singular. Then

$$U = (\mathbf{X} - \boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}) \sim \chi_p^2$$

Proof. By definition of the χ^2 -distribution we have $U = \mathbf{Z}^\top \mathbf{Z} = Z_1^2 + \dots + Z_p^2 \sim \chi_p^2$ for $\mathbf{Z} \sim N(0, \mathbf{I}_p)$. By the [Mahalanobis transformation](#) we have $\mathbf{Z} = \boldsymbol{\Sigma}^{\frac{1}{2}}(\mathbf{X} - \boldsymbol{\mu})$ from which the result follows. □

2.2 Estimation of the multivariate normal distribution

2.2.1 Univariate case

From the univariate case we recall that if $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ are independent, then the MLE estimators are:

$$\begin{aligned} \hat{\mu} &= \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2. \end{aligned}$$

But since we found that $\mathbb{E}[\sigma^2] = \frac{n-1}{n}\sigma^2$ is biased, we use instead the estimator

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

We further proved that

1. $\bar{X} \sim N(\mu, \sigma^2/n)$,
2. $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$,
3. \bar{X}, S^2 are independent (for the normal distribution),
4. $\sqrt{n} \frac{\bar{X} - \mu}{S} \sim t_{n-1}$ (student t distribution).

Our next goal is to obtain the result for the multivariate case.

2.2.2 Multivariate case

In this case, we have p -variate independent random vectors $\mathbf{X}_1, \dots, \mathbf{X}_n \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where we denote $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^\top$. These make up the columns of the $(n \times p)$ *data matrix* or *feature matrix* \mathbf{X} given as:

$$\mathbf{X} = \begin{pmatrix} X_{11} & \cdots & X_{1p} \\ \vdots & \ddots & \vdots \\ X_{n1} & \cdots & X_{np} \end{pmatrix} \begin{matrix} \text{samples} \\ \downarrow \\ \text{features} \rightarrow \end{matrix} = \begin{pmatrix} \mathbf{X}_1^\top \\ \vdots \\ \mathbf{X}_n^\top \end{pmatrix} = (\mathbf{X}_1 \cdots \mathbf{X}_n)^\top.$$

We want to estimate $\boldsymbol{\mu}, \boldsymbol{\Sigma}$. Again we denote:

$$\begin{aligned} \hat{\boldsymbol{\mu}} &= \bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i = \frac{1}{n} \mathbf{X}^\top \mathbf{1}, \\ \mathbf{S}^2 &= \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}})^\top (\mathbf{X}_i - \bar{\mathbf{X}}) = \frac{1}{n} \mathbf{X}^\top \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right) \mathbf{X}. \end{aligned}$$

The matrix $\mathbf{C} = \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top$ is called the *centering matrix* because its action is to remove the mean of a vector:

$$\mathbf{C} \mathbf{y} = \begin{pmatrix} 1 - \frac{1}{n} & \cdots & \frac{1}{n} \\ \vdots & \ddots & \vdots \\ \frac{1}{n} & \cdots & 1 - \frac{1}{n} \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix}.$$

We note that

1. \mathbf{C} is symmetric
2. \mathbf{C} is idempotent

For the estimators, we may prove:

Proposition 2. $\bar{\mathbf{X}} \sim N(\boldsymbol{\mu}, \frac{1}{n} \boldsymbol{\Sigma})$

Proof. The proof is by factoring the characteristic function. We can do this since the \mathbf{X}_i 's are independent.

$$\begin{aligned} \varphi_{\bar{\mathbf{X}}}(\mathbf{t}) &= \mathbb{E} \left[e^{i \mathbf{t}^\top \bar{\mathbf{X}}} \right] = \mathbb{E} \left[e^{i \left(\frac{\mathbf{t}}{n} \right)^\top (\mathbf{X}_1, \dots, \mathbf{X}_n)} \right] = \varphi_{\mathbf{X}_1 + \dots + \mathbf{X}_n} \left(\frac{\mathbf{t}}{n} \right) = \varphi_{\mathbf{X}_1} \left(\frac{\mathbf{t}}{n} \right) \cdots \varphi_{\mathbf{X}_n} \left(\frac{\mathbf{t}}{n} \right) \\ &= \prod_{i=1}^n e^{i \left(\frac{\mathbf{t}}{n} \right)^\top \boldsymbol{\mu} - \frac{1}{2} \left(\frac{\mathbf{t}}{n} \right)^\top \boldsymbol{\Sigma} \left(\frac{\mathbf{t}}{n} \right)} = e^{\sum_{j=1}^n i \left(\frac{\mathbf{t}}{n} \right)^\top \boldsymbol{\mu} - \frac{1}{2} \left(\frac{\mathbf{t}}{n} \right)^\top \boldsymbol{\Sigma} \left(\frac{\mathbf{t}}{n} \right)} = e^{i \mathbf{t}^\top \boldsymbol{\mu} - \frac{1}{2} \mathbf{t}^\top \frac{\boldsymbol{\Sigma}}{n} \mathbf{t}}. \end{aligned}$$

□

Proposition 3. $\mathbb{E}[\mathbf{S}] = \frac{n-1}{n}\mathbf{\Sigma}$.

Proof. This is a more straight forward computation. □

Hence we obtain an unbiased estimator as

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

2.2.3 Quadratic forms

Let $\mathbf{X} = (X_1, \dots, X_p)^\top$ be a $(p \times 1)$ random vector and $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{p \times p}$ a matrix. This gives a *quadratic form*:

$$\mathbf{X}^\top \mathbf{A} \mathbf{X} = \sum_{i,j} X_i a_{ij} X_j.$$

Theorem 12. Suppose $\mathbf{X} \sim (\boldsymbol{\mu}, \mathbf{\Sigma})$ then the *trace formula* holds:

$$\boxed{\mathbb{E}[\mathbf{X}^\top \mathbf{A} \mathbf{X}] = \text{tr}(\mathbf{A} \mathbf{\Sigma}) + \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu}} \quad (8)$$

Proof. Using that $\text{Cov}[X_i, X_j] = \mathbb{E}[X_i X_j] - \mathbb{E}[X_i] \mathbb{E}[X_j]$ we obtain the result:

$$\begin{aligned} \mathbb{E}[\mathbf{X}^\top \mathbf{A} \mathbf{X}] &= \sum_{i,j} a_{ij} \mathbb{E}[X_i X_j] = \sum_{i,j} a_{ij} (\Sigma_{ij} - \mu_i \mu_j) \\ &= \sum_{i,j} a_{ij} \Sigma_{ij} - \sum_{i,j} a_{ij} \mu_i \mu_j = \text{tr}(\mathbf{A} \mathbf{\Sigma}) - \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu}. \end{aligned}$$

□

2.2.4 Idempotent matrices

Recall that the (square) matrix \mathbf{A} is said to be idempotent if $\mathbf{A} \mathbf{A} = \mathbf{A}$. We have the following:

Proposition 4. Let \mathbf{A} be idempotent.

1. $\mathbf{I} - \mathbf{A}$ is also idempotent.
2. $\mathbf{A}(\mathbf{I} - \mathbf{A}) = (\mathbf{I} - \mathbf{A})\mathbf{A} = \mathbf{0}$.
3. The only non-singular idempotent matrix is \mathbf{I} .
4. All eigenvalues of idempotent matrices are 0 or 1.

Proof. We have:

1. $(\mathbf{I} - \mathbf{A})^2 = \mathbf{I}^2 - \mathbf{A} \mathbf{I} - \mathbf{I} \mathbf{A} + \mathbf{A}^2 = \mathbf{0}$.
2. Obvious
3. Suppose \mathbf{A} is non-singular. Then $\mathbf{I} = \mathbf{A}^{-1} \mathbf{A} = \mathbf{A}^{-1} \mathbf{A} \mathbf{A} = \mathbf{A}$.
4. $\lambda \mathbf{x} = \mathbf{A} \mathbf{x} = \mathbf{A} \mathbf{A} \mathbf{x} = \lambda^2 \mathbf{x}$.

□

Let $\mathbf{A} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^\top \in \mathbb{R}^{p \times p}$ be symmetric and idempotent. Since only non-zero eigenvalues are 1, we have $\mathbf{\Lambda} = \mathbf{I}_r$. Also $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{\Lambda}) = r$.

The following theorem will be used extensively in the sequel.

Theorem 13. $\mathbf{Z} \sim N(0, \mathbf{I}_p)$ and let \mathbf{R}, \mathbf{S} be symmetric and idempotent of rank r, s respectively. Suppose also $\mathbf{RS} = \mathbf{0}$. Then

1. $\mathbf{Z}^\top \mathbf{RZ} \sim \chi_r^2$
2. $\mathbf{Z}^\top \mathbf{RZ}$ and $\mathbf{Z}^\top \mathbf{SZ}$ are independent
3. $\frac{\mathbf{Z}^\top \mathbf{RZ}/r}{\mathbf{Z}^\top \mathbf{SZ}/s} \sim F_{r,s}$.

Proof. We have:

1. Use spectral decomposition $\mathbf{R} = \mathbf{P}\mathbf{I}_r\mathbf{P}^\top$ to find:

$$\mathbf{Z}^\top \mathbf{RZ} = (\mathbf{P}^\top \mathbf{Z})^\top \mathbf{I}_r \mathbf{P}^\top \mathbf{Z} = \mathbf{Y}^\top \mathbf{Y}.$$

Note that $\mathbf{Y} = \mathbf{P}^\top \mathbf{Z}$ is r -variate normal and compute its expectation and variance to be 0 and \mathbf{I}_r . Then finally conclude that $\mathbf{Y}^\top \mathbf{Y}^\top \mathbf{Y} = Y_1^2 + \dots + Y_r^2 \sim \chi_r^2$.

2. Compute $\text{Cov}[\mathbf{RZ}, \mathbf{SZ}] = \mathbf{RS} = \mathbf{0}$. Hence \mathbf{RZ}, \mathbf{SZ} . Then the measurable function $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{x}$ gives independence of $\mathbf{Z}^\top \mathbf{RZ}$ and $\mathbf{Z}^\top \mathbf{SZ}$.
3. By definition of Fisher distribution.

□

3 Multiple Linear Regression

3.1 Model and assumptions

We assume that we have $i = 1, \dots, n$ observations of a *response variable* Y_i depending on k *explanatory variables* x_{ij} through a linear model:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i.$$

It can be written on matrix form as:

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \\ \mathbf{Y} \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \\ \mathbf{X} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \\ \boldsymbol{\beta} \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \\ \boldsymbol{\varepsilon} \end{pmatrix}.$$

The matrix \mathbf{X} is referred to as the *design matrix*. The ε 's are *errors* and the β 's the *parameters*. We further assume:

1. \mathbf{X} is of full column rank.
2. $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$.
3. Homoscedastic: $\text{Var}[\varepsilon_i] = \sigma^2 \quad \forall i$.
4. If \mathbf{X} is random, then 2 and 3 are conditioned on \mathbf{X} .
5. Normality of errors: $\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n)$.

From the fifth assumption it follows that when \mathbf{X} is non-random we have

$$\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n).$$

We denote the estimators of $\boldsymbol{\beta}, \sigma^2$ by $\hat{\boldsymbol{\beta}}, \hat{\sigma}^2$. From these we obtain *fitted values*:

$$\hat{Y}_i = \hat{\beta}_0 + \dots + \hat{\beta}_k x_{ik} = \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}.$$

We define *residuals* by:

$$\begin{aligned} \hat{\varepsilon}_i &= Y_i - \hat{Y}_i, \\ \hat{\boldsymbol{\varepsilon}} &= \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}. \end{aligned}$$

3.2 Parameter estimation

When estimating the above parameters, there are two approaches. We may use the *least squares estimator* (LSE):

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{k+1}} \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{k+1}} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}).$$

or we may use the *maximum likelihood estimator* (MLE):

$$\hat{\beta} = \arg \max_{\beta \in \mathbb{R}^{k+1}} L(\beta), \quad L(\beta) = \prod_{i=1}^n f(Y_i).$$

It turns out that with our assumptions the result is the same. For the LSE, differentiating the sum of squares and equating to zero yields:

$$\boxed{\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}} \quad (9)$$

Having found this, we denote the *fitted values* by

$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\beta} = \underbrace{\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{\mathbf{H}} \mathbf{Y} = \mathbf{H} \mathbf{Y}.$$

The matrix \mathbf{H} is called the *prediction matrix* or *hat matrix*, and is of special interest:

Proposition 5. *For the hat matrix we have:*

1. \mathbf{H} is symmetric.
2. \mathbf{H} is idempotent.
3. $\text{rank}(\mathbf{H}) = p$.
4. Residuals can be expressed $\hat{\varepsilon} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$ with $\mathbf{I} - \mathbf{H}$ symmetric, idempotent of rank $n - p$.

Our next goal is to estimate σ^2 . More computation shows that the MLE is given by

$$\hat{\sigma}^2 = \frac{1}{n} \hat{\varepsilon}^\top \hat{\varepsilon},$$

but this is skewed as

$$\mathbb{E} [\hat{\varepsilon}^\top \hat{\varepsilon}] = \sigma^2(n - p).$$

Hence our unbiased estimator is:

$$\boxed{\hat{\sigma}^2 = \frac{1}{n - p} (\mathbf{Y} - \mathbf{X} \hat{\beta})^\top (\mathbf{Y} - \mathbf{X} \hat{\beta})}. \quad (10)$$

3.2.1 Properties of the the estimators, fitted values and residuals

We begin by remarking that $\mathbf{X} \beta$ is a linear combination of columns of \mathbf{X} and hence lies in $\text{col}(\mathbf{X})$. The same is true for $\hat{\mathbf{Y}}$.

Proposition 6. *We have:*

1. $\varepsilon \perp \hat{\mathbf{Y}}$ and $\hat{\varepsilon} \perp \text{col}(\mathbf{X})$.
2. $\sum_{i=1}^n \hat{\varepsilon}_i = 0$ and $\sum_{i=1}^n Y_i = \sum_{i=1}^n \hat{Y}_i$.
3. $\hat{\beta} \sim N(\beta, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$.
4. $\hat{\varepsilon} \sim N(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{H}))$.

5. $\frac{(n-p)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2$.

6. $\hat{\beta}, \hat{\sigma}^2$ are independent.

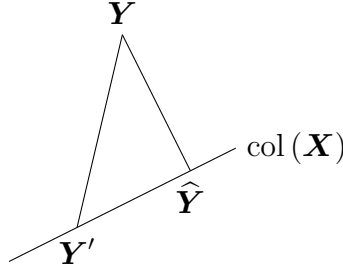


Figure 2: $\hat{\mathbf{Y}}$ is the projection onto the column space of \mathbf{X} .

Proof. **TODO: lecture 10**

□

3.2.2 Inference about β_j

Our next goal is to make confidence intervals and to perform t-tests. Recall first that the random variable T has (by definition) the *Student's t-distribution* with m degrees of freedom if it can be written as:

$$T = \frac{Z}{\sqrt{V/m}}$$

where $Z \sim N(0, 1)$, $V \sim \chi_m^2$ are independent. We may then find values for $t_{\alpha, m}$ s.t. $\mathbb{P}[T \geq t_{\alpha, m}] = \alpha$ in tables.

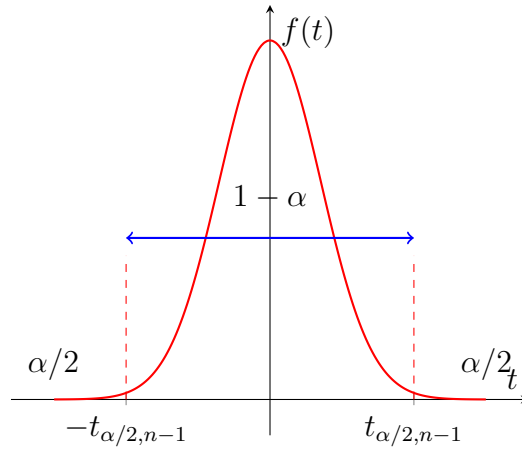


Figure 3: Two sided inference with the student-t distribution.

We have seen that $\hat{\beta} \sim N(\beta, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$. Denote $(\mathbf{X}^\top \mathbf{X})^{-1} = (e_{ij})_{i,j=1,\dots,p}$. We then have $\hat{\beta}_j \sim N(\beta_j, \sigma^2 e_{jj})$, so $\frac{\hat{\beta}_j - \beta_j}{\sqrt{e_{jj}}\sigma} \sim N(0, 1)$. Since the variance is unknown, consider the statistic:

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{e_{jj}}\hat{\sigma}} = \frac{(\hat{\beta}_j - \beta_j)/\sigma\sqrt{e_{jj}}}{\sqrt{\frac{(n-p)\hat{\sigma}^2}{\sigma^2}/(n-p)}}.$$

Recalling the properties of the estimators, we know that $\hat{\beta}, \hat{\sigma}^2$ are independent, that the numerator is $N(0, 1)$ -distributed and that $V = \frac{(n-p)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2$. Hence, we may conclude:

$$\boxed{T = \frac{\hat{\beta}_j - \beta_j}{\sqrt{e_{jj}\hat{\sigma}}} \sim t_{n-p}} \quad (11)$$

With this, we may constrict *confidence interval* by rewriting the inequalities in the expression:

$$\mathbb{P} \left[-t_{\frac{\alpha}{2}, n-p} \leq \frac{\hat{\beta}_j - \beta_j}{\sqrt{e_{jj}\hat{\sigma}}} \leq t_{\frac{\alpha}{2}, n-p} \right] = 1 - \alpha.$$

We may also perform *hypothesis testing*. Consider the following test at significance level α :

$$H_0 : \beta_j = 0, \quad H_1 : \beta_j \neq 0.$$

Under H_0 we have $T = \frac{\hat{\beta}_j - \beta_j}{\sqrt{e_{jj}\hat{\sigma}}} \sim t_{n-p}$. The *critical region* under two-sided alternative is:

$$|T| \geq t_{\frac{\alpha}{2}, n-p} \Rightarrow H_0 \text{ is rejected.}$$

TODO: R printout with explanation of columns

TODO: Does R do one or two sided hypothesis test ???

3.3 Some notes on independence

We sometimes use that plugging independent random variables through some functions result in new independent random variables. The theorem below tells us when this is okay.

Theorem 14. Suppose X, Y are independent random variables and that f, g are two measurable functions. Then $f(X), g(Y)$ are also independent.

A *measurable function* is a function s.t. the preimages of Borel sets are measurable in the given probability space. In particular, continuous functions are measurable.

3.4 Analysis of variance (ANOVA)

The following theorem forms the basis on our discussion of ANOVA.

Theorem 15. Assuming the necessary assumptions, we have the *ANOVA decomposition*:

$$\boxed{\underbrace{\sum_{i=1}^n (Y_i - \bar{Y})^2}_{\text{SST}} = \underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}_{\text{SSR}} + \underbrace{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}_{\text{SSE}}}$$

Proof. We first split the sum as:

$$\begin{aligned} \sum_{i=1}^n (Y_i - \bar{Y})^2 &= \sum_{i=1}^n (Y_i - \hat{Y}_i + \hat{Y}_i - \bar{Y})^2 \\ &= \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + 2 \sum_{i=1}^n (Y_i - \hat{Y}_i) (\hat{Y}_i - \bar{Y}). \end{aligned}$$

Using again the properties of the estimators, we find that the last sum is 0:

$$\sum_{i=1}^n (Y_i - \hat{Y}) (\hat{Y}_i - \bar{Y}) = \underbrace{\sum_{i=1}^n \overbrace{(Y_i - \hat{Y}_i)}^{\varepsilon_i} \hat{Y}_i}_{=\hat{\varepsilon}^\top \hat{\mathbf{Y}}=0} - \bar{Y} \underbrace{\sum_{i=1}^n (Y_i - \hat{Y})}_{=0 \text{ by property 2}}.$$

□

The 3 sums are called *total sum of squares*, *regression sum of squares* and *error sum of squares* respectively. This decomposition motivates the following definition.

Definition 3. The part of the total variation due to the model is called the *coefficient of determination* or the *R2-score*:

$$R^2 = \frac{\text{SSR}}{\text{SST}} \stackrel{\text{thm}}{=} 1 - \frac{\text{SSE}}{\text{SST}} \quad (12)$$

The R2-score is a measure of goodness-of-fit as it tells us how much of the variation in the data can be explained by the model. One may also prove another representation:

$$R^2 = \frac{\left(\sum_{i=1}^n (Y_i - \bar{Y})(\hat{Y}_i - \bar{Y}) \right)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2 \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}.$$

This is the square of the empirical correlation between $\mathbf{Y}, \hat{\mathbf{Y}}$.

3.4.1 Fictional model

The following discussion will examine what happens when an explanatory variable is explained by the other explanatory variables. We introduce a *fictional model* using x_{ij} as response for some fixed feature j . Wlog use feature k . The model is:

$$x_{i,k} = \alpha_0 + \alpha_1 x_{i,1} + \cdots + \alpha_{k-1} x_{i,k-1} + \delta_i.$$

As usual, we assume $\boldsymbol{\delta} = (\delta_1, \dots, \delta_n)^\top \sim N(0, \sigma^2 \mathbf{I})$. We may then estimate $\boldsymbol{\alpha} = (\alpha_0, \dots, \alpha_{k-1})^\top$ and σ^2 by the usual $\hat{\boldsymbol{\alpha}}, \hat{\sigma}^2$ to obtain fitted \hat{x}_{ik} . We find the squared empirical correlation between $x_{i,k}$ and $\hat{x}_{i,k}$ as:

$$R_k^2 = \frac{\sum_{i=1}^n (\hat{x}_{ik} - \bar{x}_k)^2}{\sum_{i=1}^n (x_{ik} - \bar{x}_k)^2} = \frac{(\sum_{i=1}^n (x_{ik} - \bar{x}_k)(\hat{x}_{ik} - \bar{x}_k))^2}{\sum_{i=1}^n (x_{ik} - \bar{x}_k)^2 \sum_{i=1}^n (\hat{x}_{ik} - \bar{x}_k)^2}.$$

We call it the coefficient of determination for x_{ik} as response. Repeating the procedure for the remaining $x_{ij}, j = 1, \dots, k-1$ we obtain R_1^2, \dots, R_k^2 . It turns out that:

$$\text{Var} [\hat{\beta}_j] = \frac{\sigma^2}{(1 - R_j^2) \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}.$$

So the more a variable is explained by the other, the higher the variance of the estimator.

3.4.2 Further expressions for the sums of squares

Recall that \mathbf{C}, \mathbf{H} are both symmetric and idempotent. For the total sum of squares, using that the centering matrix \mathbf{C} is idempotent, we obtain:

$$\text{SST} = \sum_{i=1}^n (Y_i - \bar{Y})^2 = (\mathbf{C}\mathbf{Y})^\top (\mathbf{C}\mathbf{Y}) = \mathbf{Y}^\top \mathbf{C}\mathbf{Y} = \mathbf{Y}^\top \left(\mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right) \mathbf{Y}.$$

For the residual sum of squares we also need the fact that $\mathbf{H}\mathbf{x}_i = \mathbf{x}_i$ for all columns of \mathbf{X} . This follows readily as $\mathbf{H}\mathbf{X} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} = \mathbf{X}$. From this we have $\mathbf{H}\mathbf{1} = \mathbf{1}$ as this is the first column of \mathbf{X} .

$$\begin{aligned} \text{SSR} &= \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = (\mathbf{C}\mathbf{H}\mathbf{Y})^\top (\mathbf{C}\mathbf{H}\mathbf{Y}) = \mathbf{Y}^\top \mathbf{H}\mathbf{C}\mathbf{H}\mathbf{Y} \\ &= \mathbf{Y}^\top \left(\mathbf{H} - \frac{1}{n} \mathbf{H}\mathbf{1}\mathbf{1}^\top \mathbf{H} \right) \mathbf{Y} = \mathbf{Y}^\top \left(\mathbf{H} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right) \mathbf{Y}. \end{aligned}$$

About the matrix $\mathbf{H} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top$, we note that it is symmetric, idempotent and of rank $p - 1$:

$$\text{rank} \left(\mathbf{H} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right) = \text{tr} \left(\mathbf{H} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top \right) = \text{tr}(\mathbf{H}) - \frac{1}{n} \text{tr}(\mathbf{1}\mathbf{1}^\top) = p - 1.$$

Finally, for the error sum of squares we obtain using that $\mathbf{I} - \mathbf{H}$ is symmetric and idempotent:

$$\text{SSE} = \sum_{i=1}^n (\hat{Y}_i - Y_i)^2 = \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} = ((\mathbf{I} - \mathbf{H})\mathbf{Y})^\top (\mathbf{I} - \mathbf{H})\mathbf{Y} = \dots = \mathbf{Y}^\top (\mathbf{I} - \mathbf{H})\mathbf{Y}.$$

3.5 F-test

We first consider the hypothesis test:

$$\text{H}_0 : \beta_i = 0 \quad \forall i \in \{1, \dots, k\}, \quad \text{H}_1 : \beta_i \neq 0 \text{ for at least one } i \in \{1, \dots, k\}.$$

As usual we denote the significance level by α . Our test statistic F is:

$$\boxed{F = \frac{\text{SSR}/(p-1)}{\text{SSE}/(n-p)} \sim F_{p-1, n-p}} \quad (13)$$

And as usual the test is $F \geq f_{\alpha, p-1, n-p} \Rightarrow \text{H}_0$ is rejected.

Proof. **TODO: Uke 8**

□

3.6 General F-test

We set up a much more general problem. Let $A \in \mathbb{R}^{r \times p}$, $r < p$, $\text{rank}(A) = r$, $\mathbf{d} \in \mathbb{R}^d$. We test the hypothesis:

$$\text{H}_0 : A\boldsymbol{\beta} = \mathbf{d}, \quad \text{H}_1 : A\boldsymbol{\beta} \neq \mathbf{d}.$$

Some special cases of this general setup are.

- $r = 1, d = 0, A = (0, \dots, 1, \dots, 0)$ with 1 at index i , gives the test

$$H_0 : \beta_i = 0, \quad H_1 : \beta_i \neq 0.$$

- $r = 1, d = 0, A = (0, \dots, 1, \dots, -1, \dots, 0)$ with 1 at index i and -1 at index j , gives the test

$$H_0 : \beta_i = \beta_j, \quad H_1 : \beta_i \neq \beta_j.$$

- $r = k, d = \mathbf{0} \in \mathbb{R}^k, A = (\mathbf{0}, \text{diag}(1)) \in \mathbb{R}^{k \times p}$, gives the test

$$H_0 : \beta_i = 0 \quad \forall i \in \{1, \dots, k\}, \quad H_1 : \beta_i \neq 0 \text{ for some } i \in \{1, \dots, k\}.$$

This is the F-test of the previous section.

Let \mathcal{B} be the space of β satisfying H_0 . The restricted problem is:

$$\hat{\beta}^R = \arg \min_{\beta \in \mathcal{B}} (\mathbf{Y} - \mathbf{X}\beta)^\top (\mathbf{Y} - \mathbf{X}\beta).$$

Using lagrange multipliers and a bag of tricks, we obtain:

$$\hat{\beta}^R = \hat{\beta} - (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top (\mathbf{A}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{A}^\top)^{-1} (\mathbf{A}\hat{\beta} - \mathbf{d}).$$

Denoting $\Delta = \hat{\beta} - \hat{\beta}^R$, we find:

$$\text{SSE}^R = \text{SSE} + \Delta^\top \mathbf{X}^\top \mathbf{X} \Delta$$

We claim that the under H_0 , we have

$$F = \frac{(\text{SSE}^R - \text{SSE}) / r}{\text{SSE} / (n - p)} \sim F_{r, n-p} \quad (14)$$

Proof. **TODO: Uke 8**

□

3.7 Transformations of data

Some models common in applications are:

1. $Y = \beta_0 + \beta_1 \frac{1}{x} + \varepsilon$
2. $Y = \frac{1}{\beta_0 x + \beta_1 + \varepsilon}$
3. $Y = \frac{x}{\beta_0 x + \beta_1 + x\varepsilon}$

These can also be analysed using linear regression after undergoing some *transformations*, namely:

1. $\tilde{x} = \frac{1}{x} \rightsquigarrow Y = \beta_0 + \beta_1 \tilde{x} + \varepsilon$
2. $\tilde{Y} = \frac{1}{Y} \rightsquigarrow \tilde{Y} = \beta_0 + \beta_1 x + \varepsilon$
3. $\tilde{x} = \frac{1}{x}, \tilde{Y} = \frac{1}{Y} \rightsquigarrow \tilde{Y} = \beta_0 + \beta_1 \tilde{x} + \varepsilon$

We do not always know what form the model should have, which motivates the following section.

3.7.1 Box-Cox transformation

The *Box-Cox transformation* is a power transform of the form:

$$u_\lambda(y) = \begin{cases} \frac{y^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0, \\ \ln y, & \text{if } \lambda = 0. \end{cases} \quad (15)$$

We immediately see that this requires $\mathbf{Y} > 0$ the response \mathbf{Y} . To remedy negative \mathbf{Y} , simply shift them. To apply the transform, suppose that $\tilde{Y} = u_\lambda(Y)$ is normal and choose λ which maximises the log-likelihood function. This turns out to not have analytic solution. To find the optimal λ , we therefore perform a *grid search*. In practice, we use the R function `boxcox`.

3.7.2 Variance stabilising transformation

Suppose $\mu_i = \mathbb{E}[Y_i]$ and that $\text{Var}[Y_i]$ depends on μ_i by say $\text{Var}[Y_i] = h(\mu_i)$. Our goal is to transform by $\tilde{Y} = g(Y)$ so that the variance depends less on μ . First order Taylor expansion of $g(Y)$ around μ gives:

$$g(y) \approx g(\mu) + g'(\mu)(y - \mu).$$

We find:

$$\begin{aligned} \mathbb{E}[g(Y)] &\approx g(\mu) + g'(\mu) \overbrace{(\mathbb{E}[Y] - \mu)}^{=0} = g(\mu) \\ \text{Var}[g(Y)] &\approx \text{Var}[g(\mu) + g'(\mu)(Y - \mu)] = g'(\mu)^2 \text{Var}[Y] \end{aligned}$$

Choose $g(y) = \int_c^y \frac{1}{\sqrt{h(\mu)}} d\mu$. Then $g'(y) = \frac{1}{\sqrt{h(y)}}$ so $\text{Var}[g(Y)] = 1$.

4 Model Analysis, Selection and Multiple Hypothesis Testing

4.1 Model analysis

Given a linear model, we can to some examination on the basis of visual analysis of residuals

TODO: Example QQ plots ?

TODO: not homostochastic, not independent ... transformed residuals free of trouble... TODO: long discussion of standardized / studentizes residuals ...

4.2 Model selection

Some challenges when selecting model include:

1. Unnecessary explanatory variables \rightsquigarrow overfitting
2. Missed relevant explanatory variables \rightsquigarrow underfitting
3. Quality of predictor \nleftrightarrow Quality of estimator

TODO: example from 2017 exam

Underfitting leads to biased estimation

Overfitting leads to greater variance in the estimators

4.2.1 Which model to choose

Suppose k covariates. Then we have 2^k possible models from maximal:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_k x_{ik}.$$

to minimal:

$$Y_i = \beta_0.$$

We want to arrive at a compromise between simplicity and goodness of fit. It is clear that the usual R^2 -score will not decrease by introducing additional variables, so we need to intruduce penalty for extra variables. Let $\hat{\sigma}^2$ be obtained from maximal model and SSE, p from the model of consideration.

1. *Adjusted coefficient of determination*:

$$R_{\text{adj}}^2 = 1 - \frac{SSE/(n-p-1)}{SST/(n-1)} = 1 - (1 - R^2) \frac{n-1}{n-p-1}$$

2. *Mallow's C_p parameter* ("Complexity parameter"):

$$C_p = \frac{SSE}{\hat{\sigma}^2} - n + 2p.$$

3. *Akaike information criterion*:

$$\text{AIC} = \frac{\text{SSE}/n}{\hat{\sigma}^2} + 2\frac{p}{n}.$$

4. *Bayesian information criterion*

$$\text{BIC} = \frac{\text{SSE}/n}{\hat{\sigma}^2} + \ln n \frac{p}{n}.$$

Note that for the 3 last options we want to minimise the statistic. In practice selection is done either with software or in to steps:

1. For each $j = 1, \dots, k$ choose the model with j covariates of maximal R2-score. This gives k models, so far not penalised.
2. Consider these k best models and choose “best of the best” using one of the criteria.

4.3 Multiple hypothesis testing

As before we consider the test:

$$H_0 : \beta_1 = \dots = \beta_k = 0.$$

If this is now rejected, what about the individual tests β_j ? We can test individually using t-tests the hypothesis:

$$H_0 : \beta_j = 0.$$

But if we have many parameters, $\mathbb{P}[\text{at least one type I error}]$ is large. In the worst case, for m independent hypothesis we can compute this probability to be $1 - (1 - \alpha)^m$, which at significance $\alpha = 0.05$ and $m = 5, 20$ gives a probability $> 0.22, > 0.64$ respectively. We will now consider methods to resolve this issue.

4.3.1 p -value is a random variable

We recall that given an observation t for the test statistic T we can compute the p -value, the probability of an equal or more extreme observation:

$$p(t) = \mathbb{P}_{H_0}(T \geq t).$$

Since t is the observed value of a random variable, it is clear that the p -value is also random. The rejection criteria for H_0 is:

$$H_0 \text{ is rejected} \Leftrightarrow p(t) \leq \alpha.$$

An equivalent expression for the rejection criteria is:

$$H_0 \text{ is rejected} \Leftrightarrow t \geq t' \text{ with } \mathbb{P}_{H_0}(T \geq t') = \alpha.$$

Hence we have

$$\mathbb{P}_{H_0}(p(t) \leq \alpha) = \alpha,$$

from which we conclude that under true H_0 , the p -value is uniformly distributed on $[0, 1]$.

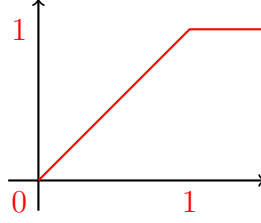


Figure 4: Distribuon of the p -value.

4.3.2 Testing m hypothesis

We begin by finding the m p -values (p_1, \dots, p_m) corresponding to each hypothesis. We then need to choose the *local significance level* α_{loc} and follow the rule that if $p\text{-value} \leq \alpha_{\text{loc}}$ we reject the corresponding H_0 . We need to choose α_{loc} such that the probability of at least one type I error is α . The outcomes of our hypothesis test can be summarised by the table: Note that we know m ,

Table 1: Multiple testing set-up

	H0 not rejected	H0 rejected	total
H_0 true	U	V	m_0
H_0 false	T	S	$m - m_0$
total	$m - R$	R	m

the number of hypothesis, and the number R of rejected hypothesis. Our goal is to control the *familywise error rate (FWER)*. This is the probability of at least one false positive finding:

$$\text{FWER} = \mathbb{P}[V \geq 1] = 1 - \mathbb{P}[V = 0].$$

For independent hypothesis we find

$$\begin{aligned} \text{FWER} &= 1 - \mathbb{P}[p_1 > \alpha_{\text{loc}}, \dots, p_m > \alpha_{\text{loc}}] \\ &= 1 - \mathbb{P}[p_1 > \alpha_{\text{loc}}] \dots \mathbb{P}[p_m > \alpha_{\text{loc}}] \\ &= 1 - (1 - \alpha_{\text{loc}})^m. \end{aligned}$$

If we for each hypothesis define the event:

$$R_j = \{\text{the } j\text{-th } H_0 \text{ is rejected but true}\} = \{p_j \leq \alpha_{\text{loc}}\}.$$

Then $\bar{R}_j = \{p_j > \alpha_{\text{loc}}\}$ is the complementary event. We may write:

$$\text{FWER} = \mathbb{P}[R_1 \cup \dots \cup R_m] = 1 - \mathbb{P}[\bar{R}_1 \cap \dots \cap \bar{R}_m].$$

We present two methods for choosing α_{loc}

1. The *Bonferrony method* uses subadditivity to bound the *FWER* by:

$$\alpha = \text{FWER} = \mathbb{P}[R_1 \cup \dots \cup R_m] \leq \sum_{j=1}^m \mathbb{P}[R_j] = m\alpha_{\text{loc}}.$$

Hence the Bonferrony method is to let:

$$\boxed{\alpha_{\text{loc}} = \frac{1}{m}\alpha} \tag{16}$$

We make the following remarks:

- (a) This gives strong control (it works under any combination of true and false hypothesis).
 - (b) To get equality we need disjoint events / perfectly negatively associated hypothesis.
 - (c) It is conservative: modelling dependencies may give lower requirement.
2. For the *Šidák method* we make the assumption that we have *independent tests*. Then we have seen that $\text{FWER} = 1 - (1 - \alpha_{\text{loc}})^m$. Hence the method is:

$$\boxed{\alpha_{\text{loc}} = 1 - (1 - \alpha)^{1/m}} \tag{17}$$

Again, we have strong control. The estimate is also exact when the tests are independent. It is conservative for positively dependent tests and liberal for negatively dependent tests. It is possible to show that the Šidák correction is always greater than the Bonferrony correction, so it is slightly less conservative; but we need to assume independence for exactness.

5 ANOVA and Design of Experiment

5.1 Analysis of variance (ANOVA)

From wikipedia; “in its simplest form, ANOVA provides a statistical test of whether two or more population means are equal, and therefore generalizes the t-test beyond two means. In other words, the ANOVA is used to test the difference between two or more means.” We will investigate this.

Suppose we have p treatments/strategies, each called a *level*, for maximising some response (health/income etc.). Suppose further that we have n_i samples of each level $i = 1, \dots, p$:

$$\begin{array}{ll} 1^{\text{st}} \text{ level} & Y_{1,1}, \dots, Y_{n_1,1} \\ 2^{\text{nd}} \text{ level} & Y_{2,1}, \dots, Y_{n_2,2} \\ \dots & \\ p^{\text{th}} \text{ level} & Y_{p,1}, \dots, Y_{n_p,p} \end{array}$$

We model

$$Y_{ij} = \mu_j + \varepsilon_{ij}, \quad 1 \leq j \leq p, \quad 1 \leq i \leq n_j.$$

We further assume that all $\varepsilon \sim N(0, \sigma^2)$ are independent. Our goal is to compare the values of μ . For instance we may want to test

$$H_0 : \mu_1 = \dots = \mu_p \quad \text{VS} \quad H_1 : \text{not all equal.}$$

We want to rewrite the problem so that we can use the tools we have developed for linear regression. It is then more convenient to use the following model. Let $\mu = \frac{1}{N} \sum_{j=1}^p n_j \mu_j$ be the overall average, where $N = n_1 + \dots + n_p$. Denote the deviation from the overall average by $\alpha_j = \mu_j - \mu$. Our model then becomes:

$$Y_{ij} = \mu + \alpha_j + \varepsilon_{ij}.$$

Note further that by definition of μ , we have $\sum_{j=1}^p n_j \alpha_j = 0$. In other words, one of the deviations can be expressed by the rest. In the following we simplify by assuming equal samples $n_1 = \dots = n_p$. Then we can express $\alpha_p = -\alpha_1 - \dots - \alpha_{p-1}$. Some more rewriting left to the reader gives the model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ where:

$$\boldsymbol{\beta} = \begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_{p-1} \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} \mathbf{1}_n & \mathbf{1}_n & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{1}_n & \mathbf{0} & \mathbf{1}_n & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{1}_n & -\mathbf{1}_n & -\mathbf{1}_n & \dots & -\mathbf{1}_n \end{pmatrix}.$$

Our hypothesis test is equivalent to

$$H_0 : \alpha_1 = \dots = \alpha_{p-1} = 0 \quad \text{VS} \quad H_1 : \text{not all 0.}$$

Hence, we may use the F -test that we developed in Section 3.5. Using the general F -test we may test many other as well.

5.1.1 Two-way ANOVA

We may also have 2 levels both present at the same time. Our model is then:

$$Y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}, \quad i = 1, \dots, r, \quad j = 1, \dots, s.$$

Adding replicates at the different combinations of i, j we have

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk}, \quad k = 1, \dots, n_{ij}.$$

Finally we may add an interaction term between the two:

$$Y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \varepsilon_{ijk}.$$

Also here we have constraints:

$$\sum_{i=1}^r \alpha_i = \sum_{j=1}^s \beta_j = \sum_{i=1}^r \gamma_{ij} = \sum_{j=1}^s \gamma_{ij} = 0.$$

5.2 Two level factorial design

We suppose we have k main factors x_1, \dots, x_k making up a model of the form:

$$\mathbf{Y} = \beta_0 + \beta_1 \mathbf{x}_1 + \dots + \beta_k \mathbf{x}_k + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2).$$

We only allow for our covariates x_i to be at *low level* and *high level*. This is modeled as $-1, +1$, and is the first of the following assumptions. The second assumption will ensure simple computation and the most precise inference about coefficients:

1. Each column has entries ± 1 .
2. The columns are orthogonal, i.e. $\mathbf{1}^T \mathbf{x}_i = \sum_{i=1}^n \mathbf{x}_{ij} = 0$ and $\mathbf{x}_i^T \mathbf{x}_j = n\delta_{ij}$.

This in particular implies that we have $\mathbf{X}^T \mathbf{X} = nI_n$. Using results from regression analysis, this significantly simplifies our estimators. We get $\hat{\beta} = \frac{1}{n} \mathbf{X}^T \mathbf{Y}$:

$$\beta_j = \frac{1}{n} \sum_{i=1}^n x_{ij} Y_j.$$

The estimate of the j 'th coefficient only depends on the j 'th covariate! Hence removal of covariates leaves the rest unchanged. We call the covariates x_1, \dots, x_k the main factors. A central question in experimentation is what happens when a factor goes from high to low level. This motivates the definition:

Definition 4. The *main effect* of main factor j is defined as:

$$\begin{aligned} \text{effect}_j &= \text{expected response at high level} - \text{expected response at low level} \\ &= \mathbb{E}[Y|x_j = 1] - \mathbb{E}[Y|x_j = -1] = 2\beta_j. \end{aligned}$$

The estimated effect is naturally by:

$$\begin{aligned} 2\hat{\beta}_j &= \frac{1}{n/2} \sum_{x_{ij}=1} Y_i - \frac{1}{n/2} \sum_{x_{ij}=-1} Y_i \\ &= \text{estimated response at high level} - \text{estimated response at low level} =: \widehat{\text{effect}}_j \end{aligned}$$

Another consequence of the second assumption is that $\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{n}I\right)$. Hence the components are independent with the same variance σ^2/n .

We also remark that we may express the regression sum of squares as:

$$\begin{aligned} \text{SSR} &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n \left(\sum_{j=0}^n x_{ij} \hat{\beta}_j - \hat{\beta}_0 \right)^2 \\ &= \sum_{i=1}^n \left(\sum_{j=1}^n x_{ij} \hat{\beta}_j + \overbrace{x_{i0}}^{=1} \hat{\beta}_0 - \hat{\beta}_0 \right)^2 = (\mathbf{X}^* \hat{\beta}^*)^T (\mathbf{X}^* \hat{\beta}^*) \\ &= \hat{\beta}^{*T} \mathbf{X}^{*T} \mathbf{X}^* \hat{\beta}^* = \hat{\beta}^{*T} nI_k \hat{\beta}^* = n \sum_{j=1}^n \hat{\beta}_j^2. \end{aligned}$$

Here \mathbf{X}^* is the feature matrix without the first column $\mathbf{1}$. The expression shows that $n\hat{\beta}_j$ may be interpreted as the amount of variability explained by factor j .

5.2.1 Modelling interactions

To go from this to a *2^k-design*, we take into account all interactions of the factors. These are modelled as products of main factors:

$$Y = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k + \beta_{1,2} x_1 x_2 + \cdots + \beta_{k-1,k} x_{k-1} x_k + \cdots + \beta_{1,2,\dots,k} x_1 \cdots x_k.$$

We extend the design matrix accordingly, and note that we still satisfy the assumptions. Hence we still have $\hat{\beta} = \frac{1}{n} \mathbf{X}^T \mathbf{Y}$.

We define the *interaction between two factors* as half the main effect of one when the other is at high level minus half the main effect when the other is at low. This encapsulates that if the change in response between levels of factor A is not independent of the level of B there is an interaction. Given the definition, it is natural to get the *estimated interaction effect* of 1 and 2 as $2\hat{\beta}_{1,2}$. To see that this is natural, rewrite to obtain:

$$2\hat{\beta}_{1,2} = \frac{1}{1/2} \underbrace{\left(\frac{1}{n/2} \sum_{\substack{x_{i1}=1 \\ x_{i2}=1}} Y_i - \frac{1}{n/2} \sum_{\substack{x_{i1}=-1 \\ x_{i2}=1}} Y_i \right)}_{\text{Estimated effect of 1 with 2 held at high}} - \frac{1}{1/2} \underbrace{\left(\frac{1}{n/2} \sum_{\substack{x_{i1}=1 \\ x_{i2}=-1}} Y_i - \frac{1}{n/2} \sum_{\substack{x_{i1}=-1 \\ x_{i2}=-1}} Y_i \right)}_{\text{Estimated effect of 1 with 2 held at low}}$$

5.2.2 Inference about effect

We have seen that $\widehat{\text{Effect}}_j \sim N(\text{Effect}_j, \frac{4\sigma^2}{n})$. To make inference about the effect/coefficients we need to make inference about σ^2 . We cannot use estimator from multiple linear regression since for MLR we have $\hat{\sigma}^2 = \frac{\text{SSE}}{n-p}$ and here $n = p$. If we don't have replicates in our experiment, to increase n , we have to resort to one of two methods.

1. We can neglect m higher order factors where it is reasonable that $\text{Effect}_j = 0$. Then we have $\widehat{\text{Effect}}_j \sim N(0, \frac{4\sigma^2}{n})$. Letting $\sigma_{\text{effect}}^2 = \frac{4\sigma^2}{n}$ we then estimate:

$$\hat{\sigma}_{\text{effect}}^2 = \frac{1}{m} \sum_{\text{negligable}} \widehat{\text{Effect}}_j^2$$

2. The idea of *Lenth's method* / *pseudo standard error* (PSE) is to compute:

$$C_1, \dots, C_m - \text{the estimated effects } \hat{A}, \hat{B}, \widehat{AB}, \dots$$

We then do the following:

- (a) Order $|C_j|$ in increasing order.
- (b) Find the median M of $|C_1|, \dots, |C_m|$ and compute $\delta_0 = 1.5M$.
- (c) Remove the effects C_j with $|C_j| \geq 2.5\delta_0$ and find the median of the rest of $|C_j|$.

Having done this we compute the PSE as

$$PSE = 1.5 \cdot \text{median} \{|C_j| : |C_j| < 2.5\delta_0\}.$$

This is then used as our estimator; $\hat{\sigma}_{PSE} = PSE$.

In the case of replications ($n \cdot 2^k$ experiment) we use results from multiple linear regression:

$$\hat{\sigma}^2 = \frac{SSE}{n - 2^k}, \quad \frac{\hat{\beta}_j}{\hat{\sigma}/\sqrt{n}} \sim t_{n-2^k}.$$

5.3 Fractional factorial design (2^{k-r} -design)

One common fractional design is the *half fraction*, where we fulfill only half of the 2^k experiments of the full design. To choose experiments perform:

1. Choose $k - 1$ main factors and consider their full 2^{k-1} design.
2. For the k -th main factor, assign the same values as one of the interactions in the above design.

Suppose for instance we are performing a 2^{4-1} design with main factors A, B, C and D. Then to the full design for A, B, C and let:

$$D = ABC.$$

We call the above the *generator of the design*. Multiplication of this equation by D gives:

$$1 = ABCD.$$

This is called the *defining relation*. From the generator or from the defining relation we may obtain the remaining *aliases* / *confounded factors* by multiplying both sides of $1 = ABCD$ by A, B, C, AB, AC, BC to obtain:

$$\begin{array}{lll} A = BCD & B = ACD & C = ABD \\ AB = CD & AC = BD & BC = AD. \end{array}$$

generator The minimal order of $LHS + RHS$ is called the *resolution* of the design. Here the resolution is IV. The greater the resolution the better. This is because at high resolution we ensure that the main factors are only confounded by higher order interactions. We note that when estimating effects, the estimator of confounded factors is the same, in our example for instance:

$$\widehat{D} = \widehat{ABC}.$$

This will encapsulate the joint effect of both D and ABC, which again shows that high resolution is good as we often neglect the higher order interactions anyways.

5.3.1 Blocking

The motivation for *blocking* is that the conditions of the experiments may change. Hence we want a clever way of deciding which experiments to do “today” etc. Suppose that there is a fixed change h to the response, i.e. that “tomorrow” we have

$$Y_i \mapsto Y_i + h.$$

If we choose experiments such that for a given factor, the number of high and low level (+ and −) are equal, the effect of this factor is unchanged by the added h . We usually consider the main factors of greatest importance. Consider for example a 2^3 experiment done over 2 days. We say it is divided in 2 *blocks*. Let block I be done with ABC at − and block II at +. We obtain the following:

Block I							Block II						
A	B	C	AB	AC	BC	ABC	A	B	C	AB	AC	BC	ABC
+	+	-	+	-	-	-	+	+	+	+	+	+	+
+	-	+	-	+	-	-	+	-	-	-	-	+	+
-	+	+	-	-	+	-	-	-	+	+	-	-	+
-	-	-	+	+	+	-	-	+	-	-	+	-	+

Figure 5: The two blocks when partitioning by the sign of ABC.

In this example our choice is good. Only the 3-factor interaction is affected by the supposed change in response since all other have equally many + and −. It would be much worse if a main factor were confounded with the block effect.

Index

- 2^k -design, 31
- absolutely continuous, 5
- Adjusted coefficient of determination, 25
- Akaike information criterion, 26
- aliases, 32
- ANOVA decomposition, 20
- Bayesian information criterion, 26
- blocking, 33
- blocks, 33
- Bonferroni method, 27
- Box-Cox transformation, 24
- centering matrix, 14
- characteristic function, 9
- coefficient of determination, 21
- conditional distribution, 5
- confidence interval, 20
- confounded factors, 32
- covariance matrix, 6
- Cramer-Wold, 10
- critical region, 20
- cumulative distribution function, 5
- data matrix, 14
- defining relation, 32
- design matrix, 17
- discrete, 5
- eigenvalue, 4
- eigenvector, 4
- empirical PCA, 8
- error sum of squares, 21
- errors, 17
- estimated interaction effect, 31
- expectation, 6
- explanatory variables, 17
- familywise error rate (FWER), 27
- feature matrix, 14
- fictional model, 21
- fitted values, 17, 18
- generator, 33
- generator of the design, 32
- grid search, 24
- half fraction, 32
- hat matrix, 18
- high level, 30
- hypothesis testing, 20
- independent, 5
- interaction between two factors, 31
- Jacobian matrix, 9
- Jordan decomposition, 4
- least squares estimator, 17
- Lenth's method, 32
- level, 29
- local significance level, 27
- low level, 30
- Mahalanobis transformation, 7, 13
- main effect, 30
- Mallow's C_p parameter, 25
- marginal distribution, 5
- maximum likelihood estimator, 18
- measurable function, 20
- moment generating function, 9
- multivariate normal, 11
- orthogonal, 4
- parameters, 17
- positive definite, 4, 7
- prediction matrix, 18
- principal component analysis, 7
- principal components, 8
- probability density function, 5
- pseudo standard error, 32
- quadratic form, 4, 15
- R²-score, 21
- random matrix, 5
- random variable, 5
- random vector, 5
- regression sum of squares, 21
- residuals, 17
- resolution, 33

response variable, 17

sample PCA, 8

standard normal vector, 11

standardize, 12

Student's t -distribution, 19

symmetric, 4, 6

total sum of squares, 21

trace, 4

trace formula, 15

transformations, 23

Šidák method, 28