

# Advanced dataanalysis and statistical modelling, Week 3

## General Linear Models - part I

Jan K. Møller, Henrik Madsen  
Poul Thyregod

DTU-Compute  
Technical University of Denmark  
DK-2800 Kgs. Lyngby  
jkmo@dtu.dk

February 2018

# Today

- The general linear model - intro
- The multivariate normal distribution
- Deviance
- Likelihood, score function and information matrix
- Estimation
- Fitted values
- Residuals
- Partitioning of variation
- Likelihood ratio tests
- The coefficient of determination

# Overview

- 1 The general linear model - intro
- 2 The multivariate normal distribution and norms
- 3 Likelihood, score function and information matrix
- 4 Estimation, fitted values and residuals
- 5 Partitioning of variation and hypothesis tests
- 6 Coefficient of determination,  $R^2$
- 7 Appendix (A proof of th. 3.5)

# The general linear model - intro

- We will use the term *classical* GLM for the General linear model to distinguish it from GLM which is used for the Generalized linear model.
- The classical GLM leads to a unique way of describing the variations of experiments with a *continuous* variable.
- The classical GLM's include
  - Regression analysis
  - Analysis of variance - ANOVA
  - Analysis of covariance - ANCOVA
- The residuals are assumed to follow a multivariate normal distribution in the classical GLM.

# The general linear model - intro

- Classical GLM's are naturally studied in the framework of the multivariate normal distribution.
- We will consider the set of  $n$  observations as a sample from a  $n$ -dimensional normal distribution.
- Under the normal distribution model, maximum-likelihood estimation of mean value parameters may be interpreted geometrically as *projection* on an appropriate subspace.
- The likelihood-ratio test statistics for model reduction may be expressed in terms of *norms* of these projections.

# General Linear Model

- A general linear model is:

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

Example (Two-way ANOVA):

	$B_1$	$B_2$	$B_3$
$A_1$	$y_{11}$	$y_{12}$	$y_{13}$
$A_2$	$y_{21}$	$y_{22}$	$y_{23}$

## Two way ANOVA (the model):

$$y_{ij} = \mu + \alpha_i + \beta_j + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim \text{i.i.d. } N(0, \sigma^2), \quad i = 1, 2, \quad j = 1, 2, 3.$$

An expanded view of this model is:

$$\begin{array}{rclclcl}
 y_{11} & = & \mu & + & \alpha_1 & & + & \beta_1 & & + & \varepsilon_{11} \\
 y_{21} & = & \mu & & & + & \alpha_2 & + & \beta_1 & & + & \varepsilon_{21} \\
 y_{12} & = & \mu & + & \alpha_1 & & & + & \beta_2 & & + & \varepsilon_{12} \\
 y_{22} & = & \mu & & & + & \alpha_2 & & + & \beta_2 & & + & \varepsilon_{22} \\
 y_{13} & = & \mu & + & \alpha_1 & & & & & + & \beta_3 & + & \varepsilon_{13} \\
 y_{23} & = & \mu & & & + & \alpha_2 & & & + & \beta_3 & + & \varepsilon_{23}
 \end{array}$$

The exact same in matrix notation (though not identifiable):

$$\underbrace{\begin{pmatrix} y_{11} \\ y_{21} \\ y_{12} \\ y_{22} \\ y_{13} \\ y_{23} \end{pmatrix}}_{\mathbf{y}} = \underbrace{\begin{pmatrix} 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 & 1 \end{pmatrix}}_{\mathbf{X}} \underbrace{\begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}}_{\boldsymbol{\beta}} + \underbrace{\begin{pmatrix} \varepsilon_{11} \\ \varepsilon_{21} \\ \varepsilon_{12} \\ \varepsilon_{22} \\ \varepsilon_{13} \\ \varepsilon_{23} \end{pmatrix}}_{\boldsymbol{\varepsilon}}$$

The default in R would be

$$\underbrace{\begin{pmatrix} y_{11} \\ y_{21} \\ y_{12} \\ y_{22} \\ y_{13} \\ y_{23} \end{pmatrix}}_{\mathbf{y}} = \underbrace{\begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 1 & 0 & 1 \end{pmatrix}}_{\mathbf{X}} \underbrace{\begin{pmatrix} \mu \\ \alpha_2 \\ \beta_2 \\ \beta_3 \end{pmatrix}}_{\boldsymbol{\beta}} + \underbrace{\begin{pmatrix} \varepsilon_{11} \\ \varepsilon_{21} \\ \varepsilon_{12} \\ \varepsilon_{22} \\ \varepsilon_{13} \\ \varepsilon_{23} \end{pmatrix}}_{\boldsymbol{\varepsilon}}$$

- $\mathbf{y}$  is the vector of all observations
- $\mathbf{X}$  is known as the *design matrix*
- $\boldsymbol{\beta}$  is the vector of parameters
- $\boldsymbol{\varepsilon}$  is a vector of independent  $N(0, \sigma^2)$  “measurement noise”
  - The vector  $\boldsymbol{\varepsilon}$  is said to follow a *multivariate normal distribution*
  - Mean vector  $\mathbf{0}$
  - Covariance matrix  $\sigma^2 \mathbf{I}$
  - Written as:  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
- $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  specifies the model, and everything can be calculated from  $\mathbf{y}$  and  $\mathbf{X}$ .



# Construction of the design matrix

In a general linear model (with both factors and covariates), it is surprisingly easy to construct the design matrix  $\mathbf{X}$ .

- For each factor: Add one column for each level, with ones in the rows where the corresponding observation is from that level, and zeros otherwise.
- For each covariate: Add one column with the measurements of the covariate.
- Remove linear dependencies (if necessary)

Example: linear regression:

$$y_i = \alpha + \beta \cdot x_i + \varepsilon$$

In matrix notation:

$$\mathbf{y} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \cdot & \cdot \\ \cdot & \cdot \\ 1 & x_n \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + \varepsilon$$

# Overview

- 1 The general linear model - intro
- 2 The multivariate normal distribution and norms**
- 3 Likelihood, score function and information matrix
- 4 Estimation, fitted values and residuals
- 5 Partitioning of variation and hypothesis tests
- 6 Coefficient of determination,  $R^2$
- 7 Appendix (A proof of th. 3.5)

# The multivariate normal distribution

Let  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^T$  be a random vector with  $Y_1, Y_2, \dots, Y_n$  independent identically distributed (iid)  $N(0, 1)$  random variables.

Note that  $E[\mathbf{Y}] = \mathbf{0}$  and the variance-covariance matrix  $\text{Var}[\mathbf{Y}] = \mathbf{I}$ .

## Definition (Multivariate normal distribution)

$\mathbf{Z}$  has an  $k$ -dimensional multivariate normal distribution if  $\mathbf{Z}$  has the same distribution as  $\mathbf{A}\mathbf{Y} + \mathbf{b}$  for some  $n$ , some  $k \times n$  matrix  $\mathbf{A}$ , and some  $k$  vector  $\mathbf{b}$ . We indicate the multivariate normal distribution by writing  $\mathbf{Z} \sim N(\mathbf{b}, \mathbf{A}\mathbf{A}^T)$ .

Since  $\mathbf{A}$  and  $\mathbf{b}$  are fixed, we have  $E[\mathbf{Z}] = \mathbf{b}$  and  $\text{Var}[\mathbf{Z}] = \mathbf{A}\mathbf{A}^T$ .

# The multivariate normal distribution

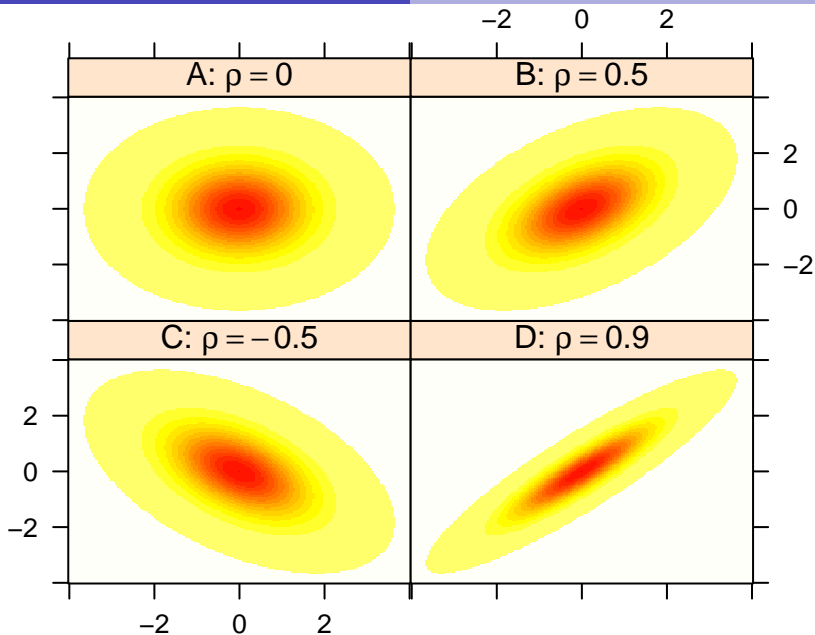
Let us assume that the variance-covariance matrix is known apart from a constant factor,  $\sigma^2$ , i.e.  $\text{Var}[\mathbf{Z}] = \sigma^2 \mathbf{\Sigma}$ .

The density for the  $k$ -dimensional random vector  $\mathbf{Z}$  with mean  $\boldsymbol{\mu}$  and covariance  $\sigma^2 \mathbf{\Sigma}$  is:

$$f_{\mathbf{Z}}(\mathbf{z}) = \frac{1}{(2\pi)^{k/2} \sigma^k \sqrt{\det \mathbf{\Sigma}}} \exp \left[ -\frac{1}{2\sigma^2} (\mathbf{z} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{z} - \boldsymbol{\mu}) \right]$$

where  $\mathbf{\Sigma}$  is seen to be (a) symmetric and (b) positive semi-definite.

We write  $\mathbf{Z} \sim N_k(\boldsymbol{\mu}, \sigma^2 \mathbf{\Sigma})$ .



# The normal density as a statistical model

Consider now the  $n$  observations  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_n)^T$ , and assume that a statistical model is

$$\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \sigma^2 \boldsymbol{\Sigma}) \text{ for } \mathbf{y} \in \mathbb{R}^n$$

The variance-covariance matrix for the observations is called the *dispersion matrix*, denoted  $D[\mathbf{Y}]$ , i.e. the dispersion matrix for  $\mathbf{Y}$  is

$$D[\mathbf{Y}] = \sigma^2 \boldsymbol{\Sigma}$$

# Inner product and norm

## Definition (Inner product and norm)

The bilinear form

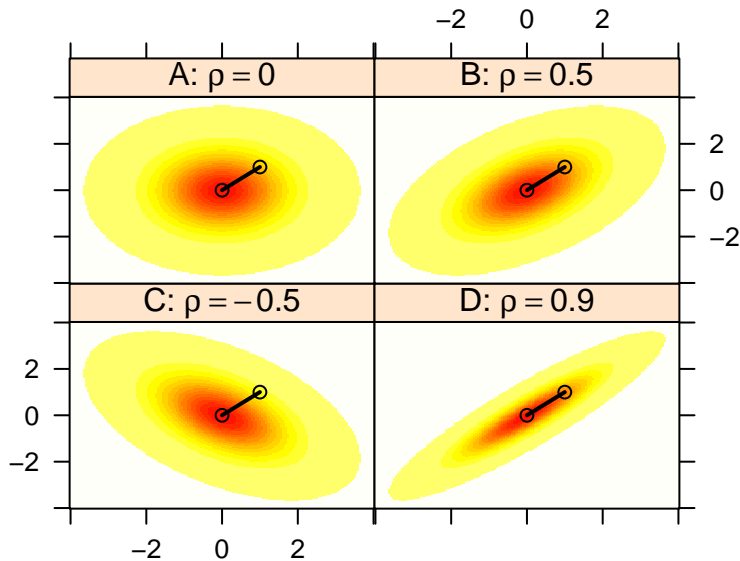
$$\delta_{\Sigma}(\mathbf{y}_1, \mathbf{y}_2) = \mathbf{y}_1^T \Sigma^{-1} \mathbf{y}_2$$

defines an *inner product* in  $\mathbb{R}^n$ . Corresponding to this inner product we can define *orthogonality*, which is obtained when the inner product is zero.

A *norm* is defined by

$$\|\mathbf{y}\|_{\Sigma} = \sqrt{\delta_{\Sigma}(\mathbf{y}, \mathbf{y})}.$$

## Order the vectors according to length (norm)





# Deviance for normal distributed variables

## Definition (Deviance for normal distributed variables)

Let us introduce the notation

$$D(\mathbf{y}; \boldsymbol{\mu}) = \delta_{\Sigma}(\mathbf{y} - \boldsymbol{\mu}, \mathbf{y} - \boldsymbol{\mu}) = (\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

to denote the quadratic norm of the vector  $(\mathbf{y} - \boldsymbol{\mu})$  corresponding to the inner product defined by  $\boldsymbol{\Sigma}^{-1}$ .

For a normal distribution with  $\boldsymbol{\Sigma} = \mathbf{I}$ , the deviance is just the Residual Sum of Squares (RSS).

## Deviance for normal distributed variables

Using this notation the normal density is expressed as a density defined on any finite dimensional vector space equipped with the inner product,  $\delta_\Sigma$ :

$$f(\mathbf{y}; \boldsymbol{\mu}, \sigma^2) = \frac{1}{(\sqrt{2\pi})^n \sigma^n \sqrt{\det(\boldsymbol{\Sigma})}} \exp \left[ -\frac{1}{2\sigma^2} D(\mathbf{y}; \boldsymbol{\mu}) \right].$$

# Overview

- 1 The general linear model - intro
- 2 The multivariate normal distribution and norms
- 3 Likelihood, score function and information matrix**
- 4 Estimation, fitted values and residuals
- 5 Partitioning of variation and hypothesis tests
- 6 Coefficient of determination,  $R^2$
- 7 Appendix (A proof of th. 3.5)

# The likelihood and log-likelihood function

- The likelihood function is:

$$L(\boldsymbol{\mu}, \sigma^2; \mathbf{y}) = \frac{1}{(\sqrt{2\pi})^n \sigma^n \sqrt{\det(\boldsymbol{\Sigma})}} \exp \left[ -\frac{1}{2\sigma^2} D(\mathbf{y}; \boldsymbol{\mu}) \right]$$

- The log-likelihood function is (apart from an additive constant):

$$\begin{aligned} \ell(\boldsymbol{\mu}, \sigma^2; \mathbf{y}) &= -(n/2) \log(\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \\ &= -(n/2) \log(\sigma^2) - \frac{1}{2\sigma^2} D(\mathbf{y}; \boldsymbol{\mu}). \end{aligned}$$

## The score function, observed - and expected information for $\mu$

- The score function wrt.  $\mu$  is

$$\frac{\partial}{\partial \mu} \ell(\mu, \sigma^2; \mathbf{y}) = \frac{1}{\sigma^2} [\Sigma^{-1} \mathbf{y} - \Sigma^{-1} \mu] = \frac{1}{\sigma^2} \Sigma^{-1} (\mathbf{y} - \mu)$$

- The observed information (wrt.  $\mu$ ) is

$$j(\mu; \mathbf{y}) = \frac{1}{\sigma^2} \Sigma^{-1}.$$

- It is seen that the observed information does not depend on the observations  $\mathbf{y}$ . Hence the expected information is

$$i(\mu) = \frac{1}{\sigma^2} \Sigma^{-1}.$$

# The general linear model

In the case of a normal density the observation  $Y_i$  is most often written as

$$Y_i = \mu_i + \epsilon_i$$

which for all  $n$  observations  $(Y_1, Y_2, \dots, Y_n)$  can be written on the matrix form

$$\mathbf{Y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

where

$$\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \sigma^2 \boldsymbol{\Sigma}) \text{ for } \mathbf{y} \in \mathbb{R}^n$$

In many applications we will assume that  $\boldsymbol{\Sigma} = \mathbf{I}$ .

# General Linear Models

- In the *linear model* it is assumed that  $\mu$  belongs to a linear (or affine) subspace  $\Omega_0$  of  $\mathbb{R}^n$ .
- The *full model* is a model with  $\Omega_{full} = \mathbb{R}^n$  and hence each observation fits the model perfectly, i.e.  $\hat{\mu} = y$ .
- The most restricted model is the *null model* with  $\Omega_{null} = \mathbb{R}$ . It only describes the variations of the observations by a common mean value for all observations.
- In practice, one often starts with formulating a rather comprehensive model with  $\Omega = \mathbb{R}^k$ , where  $k < n$ . We will call such a model a *sufficient model*.

# The General Linear Model

## Definition (The general linear model)

Assume that  $Y_1, Y_2, \dots, Y_n$  is normally distributed as described before. A *general linear model* for  $Y_1, Y_2, \dots, Y_n$  is a model where an affine hypothesis is formulated for  $\mu$ . The hypothesis is of the form

$$\mathcal{H}_0 : \mu - \mu_0 \in \Omega_0,$$

where  $\Omega_0$  is a linear subspace of  $\mathbb{R}^n$  of dimension  $k$ , and where  $\mu_0$  denotes a vector of *known offset values*.

## Definition (Dimension of general linear model)

The dimension of the subspace  $\Omega_0$  for the linear model is the *dimension of the model*.



# The design matrix

## Definition (Design matrix for classical GLM)

Assume that the linear subspace  $\Omega_0 = \text{span}\{x_1, \dots, x_k\}$ , i.e. the subspace is spanned by  $k$  vectors ( $k < n$ ).

Consider a general linear model where the hypothesis can be written as

$$\mathcal{H}_0 : \mu - \mu_0 = \mathbf{X}\boldsymbol{\beta} \text{ with } \boldsymbol{\beta} \in \mathbb{R}^k,$$

where  $\mathbf{X}$  has full rank. The  $n \times k$  matrix  $\mathbf{X}$  of known deterministic coefficients is called the *design matrix*.

The  $i^{th}$  row of the design matrix is given by the *model vector*

$$\mathbf{x}_i^T = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ik} \end{pmatrix}^T,$$

for the  $i^{th}$  observation.

# Overview

- 1 The general linear model - intro
- 2 The multivariate normal distribution and norms
- 3 Likelihood, score function and information matrix
- 4 Estimation, fitted values and residuals**
- 5 Partitioning of variation and hypothesis tests
- 6 Coefficient of determination,  $R^2$
- 7 Appendix (A proof of th. 3.5)

# Estimation of mean value parameters

Under the hypothesis

$$\mathcal{H}_0 : \boldsymbol{\mu} \in \Omega_0 ,$$

the maximum likelihood estimate for the set  $\boldsymbol{\mu}$  is found as the orthogonal projection (with respect to  $\delta_{\Sigma}$ ),  $p_0(\mathbf{y})$  of  $\mathbf{y}$  onto the linear subspace  $\Omega_0$ .

## Theorem (ML estimates of mean value parameters)

*For hypothesis of the form*

$$\mathcal{H}_0 : \boldsymbol{\mu}(\boldsymbol{\beta}) = \mathbf{X}\boldsymbol{\beta}$$

*the maximum likelihood estimated for  $\boldsymbol{\beta}$  is found as a solution to the normal equation*

$$\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{y} = \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} \hat{\boldsymbol{\beta}}.$$

*If  $\mathbf{X}$  has full rank, the solution is uniquely given by*

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{y}$$

# Properties of the ML estimator

## Theorem (Properties of the ML estimator)

*For the ML estimator we have*

$$\hat{\beta} \sim N_k(\beta, \sigma^2 (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1})$$

## Unknown $\Sigma$

Notice that it has been assumed that  $\Sigma$  is known. If  $\Sigma$  is unknown different possibilities exist

- The relaxation algorithm described in Madsen (2008) <sup>a</sup>.
- Likelihood based method, either direct or by restricted maximum likelihood (eg. `glms` in the `nlme` package of R)

---

<sup>a</sup>Madsen, H. (2008) Time Series Analysis. Chapman, Hall

# Expectation and variance of MLE

The ML estimate of  $\beta$  is central

$$\begin{aligned} E[\hat{\beta}] &= E[(\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} \mathbf{Y}] \\ &= (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} E[\mathbf{X}\beta + \epsilon] \\ &= (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} \mathbf{X}\beta = \beta \end{aligned}$$

with variance

$$\begin{aligned} V[\hat{\beta}] &= V[(\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} \mathbf{Y}] \\ &= (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} V[\mathbf{X}\beta + \epsilon] \Sigma^{-1} \mathbf{X} (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} \Sigma \Sigma^{-1} \mathbf{X} (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \\ &= \sigma^2 (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \end{aligned}$$

# Fitted values

## Fitted – or predicted – values

The *fitted* values  $\hat{\mu} = \mathbf{X}\hat{\beta}$  is found as the projection of  $\mathbf{y}$  (denoted  $p_0(\mathbf{y})$ ) on to the subspace  $\Omega_0$  spanned by  $\mathbf{X}$ , and  $\hat{\beta}$  denotes the local coordinates for the projection.

## Definition (Projection matrix)

A matrix  $\mathbf{H}$  is a *projection matrix* if and only if

- (a)  $\mathbf{H}^T = \mathbf{H}$  and
- (b)  $\mathbf{H}^2 = \mathbf{H}$ , i.e. the matrix is *idempotent*.

# The hat matrix

- Consider the model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}; \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

- The matrix

$$\mathbf{H} = \mathbf{X}[\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T$$

is a projection matrix.

- The projection matrix provides the predicted values  $\hat{\boldsymbol{\mu}}$ , since

$$\hat{\boldsymbol{\mu}} = p_0(\mathbf{y}) = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{H}\mathbf{y}$$

- Furthermore

- $\mathbf{I} - \mathbf{H}$  is also a projection matrix

- $\text{Tr}(\mathbf{H}) = p$ ,

- $\text{Tr}(\mathbf{I} - \mathbf{H}) = n - p$

# Projections - PROOFS

PROOF:

$$1: H^2 = X(X^T X)^{-1} X^T X (X^T X)^{-1} X^T$$

$$= X(X^T X)^{-1} X^T = H$$

$$H^T = (X(X^T X)^{-1} X^T)^T = ((X^T X)^{-1} X^T)^T X^T$$

$$= X((X^T X)^{-1})^T X^T = X(X^T X)^{-1} X^T = H$$

$$2: (I - H)^2 = I + H^2 - IH - HI = I - H$$

$$(I - H)^T = I^T - H^T = I - H$$

$$3^1: Tr(H) = Tr(X(X^T X)^{-1} X^T)$$

$$= Tr((X^T X)^{-1} (X^T X)) = Tr(I_p) = p$$

$$4: Tr(I - H) = Tr(I) - Tr(H) = n - p$$

---

<sup>1</sup>We use the relation  $Tr(ABC) = Tr(BCA) = Tr(CAB)$



# The hat matrix

- It follows that the predicted values are normally distributed with

$$D[\mathbf{X}\hat{\boldsymbol{\beta}}] = \sigma^2 \mathbf{X}[\mathbf{X}^T \mathbf{X}]^{-1} \mathbf{X}^T = \sigma^2 \mathbf{H}$$

- The matrix  $\mathbf{H}$  is often termed the *hat matrix* since it transforms the observations  $\mathbf{y}$  to their predicted values symbolized by a "hat" on the  $\mu$ 's.
- The observed *residuals* are

$$\mathbf{r} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{I} - \mathbf{H})\mathbf{y}$$

- The distribution of the residuals is also normal with

$$D[\mathbf{r}(\mathbf{Y})] = \sigma^2(\mathbf{I} - \mathbf{H})$$

# Residuals

## Orthogonality

The maximum likelihood estimate for  $\beta$  is found as the value of  $\beta$  which minimizes the *distance*  $\|\mathbf{y} - \mathbf{X}\beta\|$ .

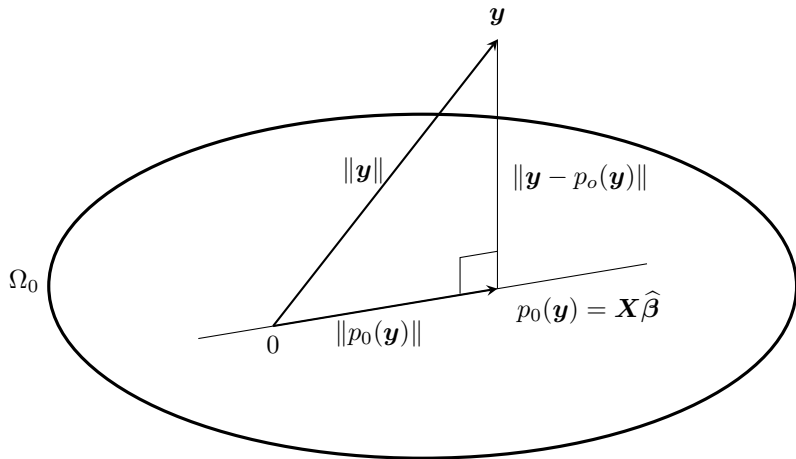
The normal equations show that (still assuming  $\Sigma = I$ )

$$\begin{aligned}\delta(\mathbf{y} - \hat{\mathbf{y}}, \hat{\mathbf{y}}) &= (\mathbf{y} - \hat{\mathbf{y}})^T \hat{\mathbf{y}} \\ &= \mathbf{y}^T (\mathbf{I} - \mathbf{H})^T \mathbf{H} \mathbf{y} \\ &= \mathbf{y}^T (\mathbf{H} - \mathbf{H}) \mathbf{y} = 0\end{aligned}$$

i.e. the *residuals* are orthogonal (with respect to  $\Sigma^{-1}$ ) to the subspace  $\Omega_0$ .

The residuals are thus orthogonal to the fitted – or predicted – values.

# Residuals



**Figure:** Orthogonality between the residual  $(\mathbf{y} - \mathbf{X}\hat{\beta})$  and the vector  $\mathbf{X}\hat{\beta}$ .

# Residuals

- The residuals  $\mathbf{r} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$  are normally distributed with

$$\mathbf{D}[\mathbf{r}] = \sigma^2(\mathbf{I} - \mathbf{H})$$

- The individual residuals do not have the same variance.
- The residuals are thus belonging to a subspace of dimension  $n - k$ , which is orthogonal to  $\Omega_0$ .
- It may be shown that the distribution of the residuals  $\mathbf{r}$  is independent of the fitted values  $\mathbf{X}\hat{\boldsymbol{\beta}}$  (follows from orthogonality and normality).

## Generalization

The results above apply to the model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}; \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

In order to generalise to the model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}; \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \boldsymbol{\Sigma})$$

Consider the “square root” ( $\mathbf{T}$ ) of  $\boldsymbol{\Sigma} = \mathbf{T}\mathbf{T}^T$ , and rewrite the model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{T}\mathbf{e}; \quad \mathbf{e} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

to ( $\mathbf{Z} = \mathbf{T}^{-1}\mathbf{Y}$ )

$$\mathbf{Z} = \mathbf{T}^{-1}\mathbf{X}\boldsymbol{\beta} + \mathbf{e}; \quad \mathbf{e} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

then the results apply to  $\mathbf{Z}$  and back transformation will give the results in the original domain.

## Example<sup>2</sup>

Two items  $A$  and  $B$  are weighed on a balance, first separately then together, giving the observations  $y_1, y_2, y_3$ , and the model

$$Y_1 = \beta_A + \epsilon_1$$

$$Y_2 = \beta_B + \epsilon_2$$

$$Y_3 = \beta_A + \beta_B + \epsilon_3$$

with  $\epsilon_i \sim N(0, \sigma^2)$ .

---

<sup>2</sup>from Bingham and Fry (2010)

# Example

Or in matrix notation

$$\begin{aligned}\mathbf{Y} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \beta_A \\ \beta_B \end{bmatrix} + \epsilon \\ &= \mathbf{X}\boldsymbol{\beta} + \epsilon\end{aligned}$$

with  $\epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$ . Hence

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \frac{1}{3} \begin{bmatrix} 2 & -1 & 1 \\ -1 & 2 & 1 \end{bmatrix} \mathbf{y}\end{aligned}$$

# Example

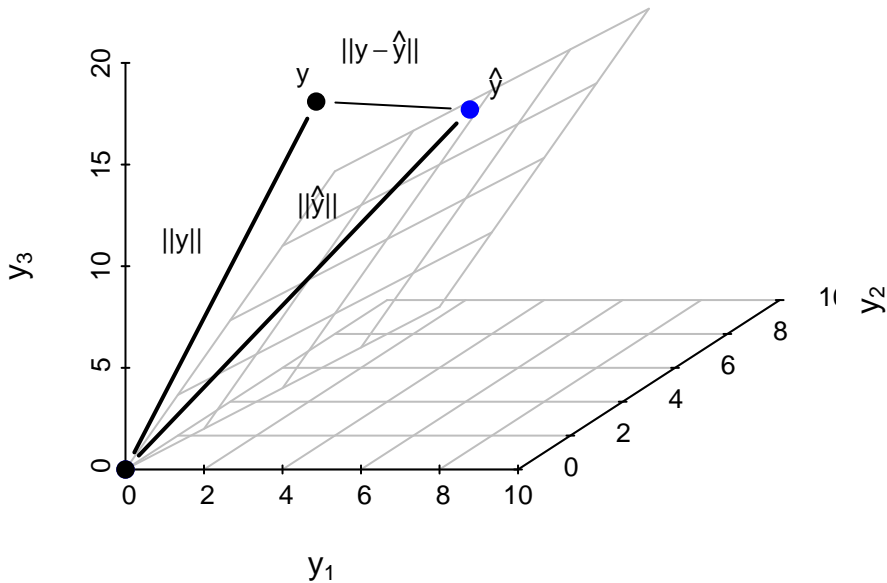
And

$$\begin{aligned}\hat{\mathbf{y}} &= \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \frac{1}{3} \begin{bmatrix} 2 & -1 & 1 \\ -1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix} \mathbf{y} \\ &= \mathbf{H} \mathbf{y}\end{aligned}$$

The projection  $\mathbf{H}$  defines a 2-dimensional surface in  $\mathbb{R}^3$ .

Exercise: Formulate a model  $y_i = x_{1,i}\beta_1 + x_{2,i}\beta_2$  such that the test  $\beta_A = \beta_B$  is equivalent to testing  $\beta_2 = 0$





# Overview

- 1 The general linear model - intro
- 2 The multivariate normal distribution and norms
- 3 Likelihood, score function and information matrix
- 4 Estimation, fitted values and residuals
- 5 Partitioning of variation and hypothesis tests**
  - Estimation of the residual variance  $\sigma^2$
- 6 Coefficient of determination,  $R^2$
- 7 Appendix (A proof of th. 3.5)

# Cochran's theorem

## Theorem (Cochran's theorem)

Suppose that  $\mathbf{Y} \sim N_n(\mathbf{0}, \mathbf{I}_n)$  (i.e. standard multivariate Gaussian random variable)

$$\mathbf{Y}^T \mathbf{Y} = \mathbf{Y}^T \mathbf{H}_1 \mathbf{Y} + \mathbf{Y}^T \mathbf{H}_2 \mathbf{Y} + \cdots + \mathbf{Y}^T \mathbf{H}_k \mathbf{Y}$$

where  $\mathbf{H}_i$  is a symmetric  $n \times n$  matrix with rank  $n_i$ ,  $i = 1, 2, \dots, k$ .  
Then any one of the following conditions implies the other two:

- i The ranks of the  $\mathbf{H}_i$  adds to  $n$ , i.e.  $\sum_{i=1}^k n_i = n$
- ii Each quadratic form  $\mathbf{Y}^T \mathbf{H}_i \mathbf{Y} \sim \chi_{n_i}^2$  (thus the  $\mathbf{H}_i$  are positive semidefinite)
- iii All the quadratic forms  $\mathbf{Y}^T \mathbf{H}_i \mathbf{Y}$  are independent (necessary and sufficient condition).

# Partitioning of variation

## Partitioning of the variation

$$\begin{aligned} D(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}) &= D(\mathbf{y}; \mathbf{X}\hat{\boldsymbol{\beta}}) + D(\mathbf{X}\hat{\boldsymbol{\beta}}; \mathbf{X}\boldsymbol{\beta}) \\ &= (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + \\ &\quad (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \\ &\geq (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \end{aligned}$$

# Partitioning of variation

## $\chi^2$ -distribution of individual contributions

Under  $\mathcal{H}_0$  it follows from the normal distribution of  $\mathbf{Y}$  that

$$D(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \sim \sigma^2 \chi_n^2$$

Furthermore, it follows from the normal distribution of  $\mathbf{r}$  and of  $\hat{\boldsymbol{\beta}}$  that

$$D(\mathbf{y}; \mathbf{X}\hat{\boldsymbol{\beta}}) = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \sim \sigma^2 \chi_{n-k}^2$$

$$D(\mathbf{X}\hat{\boldsymbol{\beta}}; \mathbf{X}\boldsymbol{\beta}) = (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sim \sigma^2 \chi_k^2$$

moreover, the independence of  $\mathbf{r}$  and  $\mathbf{X}\hat{\boldsymbol{\beta}}$  implies that  $D(\mathbf{y}; \mathbf{X}\hat{\boldsymbol{\beta}})$  and  $D(\mathbf{X}\hat{\boldsymbol{\beta}}; \mathbf{X}\boldsymbol{\beta})$  are independent.

Thus, the  $\sigma^2 \chi_n^2$ -distribution on the left side is partitioned into two independent  $\chi^2$  distributed variables with  $n - k$  and  $k$  degrees of freedom, respectively.

# Estimation of the residual variance $\sigma^2$

## Theorem (Estimation of the variance)

*Under the hypothesis*

$$\mathcal{H}_0 : \boldsymbol{\mu}(\boldsymbol{\beta}) = \mathbf{X}\boldsymbol{\beta}$$

*the maximum marginal likelihood estimator for the variance  $\sigma^2$  is*

$$\hat{\sigma}^2 = \frac{D(\mathbf{y}; \mathbf{X}\hat{\boldsymbol{\beta}})}{n - k} = \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{n - k}$$

*Under the hypothesis,  $\hat{\sigma}^2 \sim \sigma^2 \chi_f^2 / f$  with  $f = n - k$ .*

## Likelihood ratio tests

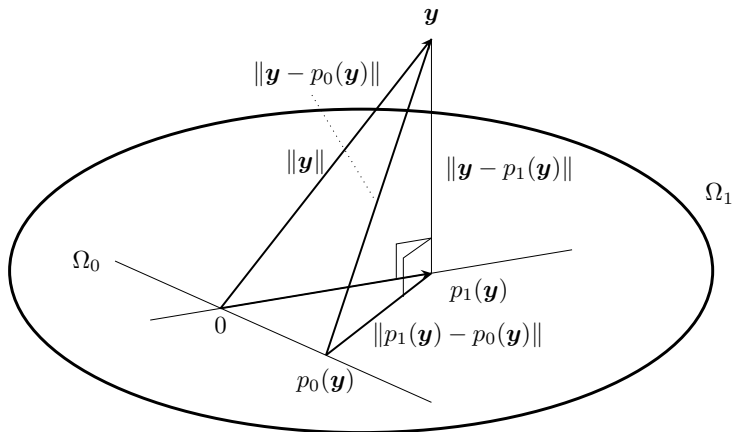
- In the classical GLM case the exact distribution of the likelihood ratio test statistic may be derived.
- Consider the following model for the data  $\mathbf{Y} \sim N_n(\boldsymbol{\mu}, \sigma^2 \boldsymbol{\Sigma})$ .
- Let us assume that we have the sufficient model

$$\mathcal{H}_1 : \boldsymbol{\mu} \in \Omega_1 \subset \mathbb{R}^n$$

with  $\dim(\Omega_1) = m_1$ .

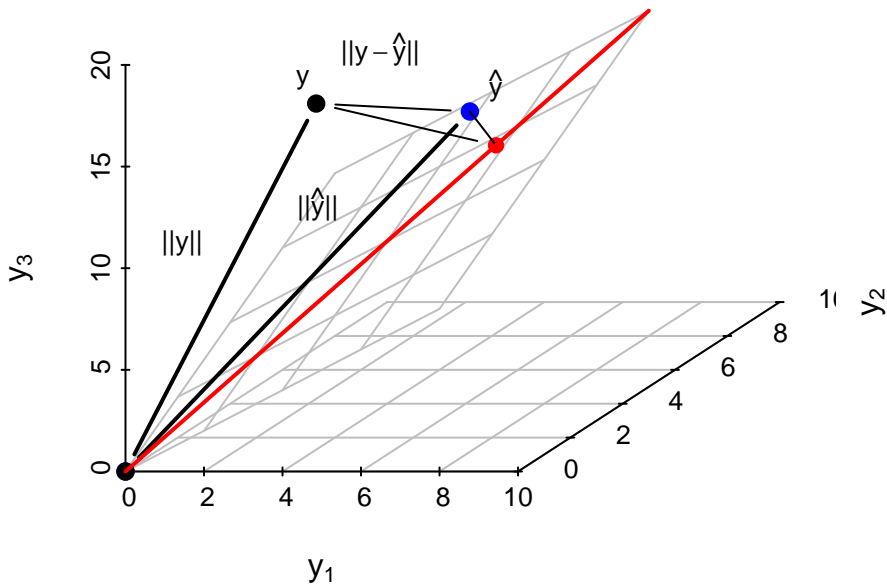
- Now we want to test whether the model may be reduced to a model where  $\boldsymbol{\mu}$  is restricted to some subspace of  $\Omega_1$ , and hence we introduce  $\Omega_0 \subset \Omega_1$  as a linear (affine) subspace with  $\dim(\Omega_0) = m_0$ .

# Model reduction



**Figure:** Model reduction. The partitioning of the deviance corresponding to a test of the hypothesis  $\mathcal{H}_0 : \mu \in \Omega_0$  under the assumption of  $\mathcal{H}_1 : \mu \in \Omega_1$ .





# Test for model reduction

## Theorem (A test for model reduction)

*The likelihood ratio test statistic for testing*

$$\mathcal{H}_0 : \boldsymbol{\mu} \in \Omega_0 \text{ against the alternative } \mathcal{H}_1 : \boldsymbol{\mu} \in \Omega_1 \setminus \Omega_0$$

*is a monotone function of*

$$F(\mathbf{y}) = \frac{D(p_1(\mathbf{y}); p_0(\mathbf{y})) / (m_1 - m_0)}{D(\mathbf{y}; p_1(\mathbf{y})) / (n - m_1)}$$

*where  $p_1(\mathbf{y})$  and  $p_0(\mathbf{y})$  denote the projection of  $\mathbf{y}$  on  $\Omega_1$  and  $\Omega_0$ , respectively. Under  $\mathcal{H}_0$  we have*

$$F \sim F(m_1 - m_0, n - m_1)$$

*i.e. large values of  $F$  reflects a conflict between the data and  $\mathcal{H}_0$ , and hence lead to rejection of  $\mathcal{H}_0$ . The  $p$ -value of the test is found as*

*$p = P[F(m_1 - m_0, n - m_1) \geq F_{obs}]$ , where  $F_{obs}$  is the observed value of  $F$  given the data.*

## Test for model reduction

- The partitioning of the variation is presented in a Deviance table (or an *ANalysis Of VAriance table*, ANOVA).
- The table reflects the partitioning in the test for *model reduction*.
- The deviance between the variation of the model from the hypothesis is measured using the deviance of the observations from the model as a reference.
- Under  $\mathcal{H}_0$  they are both  $\chi^2$  distributed, orthogonal and thus independent.
- This means that the ratio is  $F$  distributed.
- If the test quantity is large this shows evidence against the model reduction tested using  $\mathcal{H}_0$ .

# Deviance table

Source	$f$	Deviance	Test statistic, $F$
Model versus hypothesis	$m_1 - m_0$	$\ p_1(\mathbf{y}) - p_0(\mathbf{y})\ ^2$	$\frac{\ p_1(\mathbf{y}) - p_0(\mathbf{y})\ ^2 / (m_1 - m_0)}{\ \mathbf{y} - p_1(\mathbf{y})\ ^2 / (n - m_1)}$
Residual under model	$n - m_1$	$\ \mathbf{y} - p_1(\mathbf{y})\ ^2$	
Residual under hypothesis	$n - m_0$	$\ \mathbf{y} - p_0(\mathbf{y})\ ^2$	

**Table:** Deviance table corresponding to a test for model reduction as specified by  $\mathcal{H}_0$ . For  $\Sigma = \mathbf{I}$  this corresponds to an analysis of variance table, and then 'Deviance' is equal to the 'Sum of Squared deviations (SS)'

# Test for model reduction

## The test is a conditional test

It should be noted that the test has been derived as a *conditional test*. It is a test for the hypothesis  $\mathcal{H}_0 : \boldsymbol{\mu} \in \Omega_0$  under the assumption that  $\mathcal{H}_1 : \boldsymbol{\mu} \in \Omega_1$  is true. The test does in no way assess whether  $\mathcal{H}_1$  is in agreement with the data. On the contrary in the test the residual variation under  $\mathcal{H}_1$  is used to estimate  $\sigma^2$ , i.e. to assess  $D(\mathbf{y}; p_1(\mathbf{y}))$ .

## The test does not depend on the particular parametrization of the hypotheses

Note that the test does only depend on the two sub-spaces  $\Omega_1$  and  $\Omega_0$ , but not on how the subspaces have been parametrized (the particular choice of basis, i.e. the design matrix). Therefore it is sometimes said that the test is *coordinate free*.

## Initial test for model 'sufficiency'

- In practice, one often starts with formulating a rather comprehensive model, a *sufficient model*, and then tests whether the model may be reduced to the *null model* with  $\Omega_{null} = \mathbb{R}$ , i.e.  $\dim \Omega_{null} = 1$ .

- The hypotheses are

$$\mathcal{H}_{null} : \boldsymbol{\mu} \in \mathbb{R}$$

$$\mathcal{H}_1 : \boldsymbol{\mu} \in \Omega_1 \setminus \mathbb{R}.$$

where  $\dim \Omega_1 = k$ .

- The hypothesis is a hypothesis of "Total homogeneity", namely that all observations are satisfactorily represented by their common mean.

# Deviance table

Source	$f$	Deviance	Test statistic, $F$
Model $\mathcal{H}_{null}$	$k - 1$	$\ p_1(\mathbf{y}) - p_{null}(\mathbf{y})\ ^2$	$\frac{\ p_1(\mathbf{y}) - p_{null}(\mathbf{y})\ ^2 / (k - 1)}{\ \mathbf{y} - p_1(\mathbf{y})\ ^2 / (n - k)}$
Residual under $\mathcal{H}_1$	$n - k$	$\ \mathbf{y} - p_1(\mathbf{y})\ ^2$	
Total	$n - 1$	$\ \mathbf{y} - p_{null}(\mathbf{y})\ ^2$	

**Table:** Deviance table corresponding to the test for model reduction to the null model.

Under  $\mathcal{H}_{null}$ ,  $F \sim F(k - 1, n - k)$ , and hence large values of  $F$  would indicate rejection of the hypothesis  $\mathcal{H}_{null}$ . The  $p$ -value of the test is  $p = P[F(k - 1, n - k) \geq F_{obs}]$ .

# Overview

- 1 The general linear model - intro
- 2 The multivariate normal distribution and norms
- 3 Likelihood, score function and information matrix
- 4 Estimation, fitted values and residuals
- 5 Partitioning of variation and hypothesis tests
- 6 Coefficient of determination,  $R^2$**
- 7 Appendix (A proof of th. 3.5)



# Coefficient of determination, $R^2$

- The *coefficient of determination*,  $R^2$ , is defined as

$$R^2 = \frac{D(p_1(\mathbf{y}); p_{null}(\mathbf{y}))}{D(\mathbf{y}; p_{null}(\mathbf{y}))} = 1 - \frac{D(\mathbf{y}; p_1(\mathbf{y}))}{D(\mathbf{y}; p_{null}(\mathbf{y}))}, \quad 0 \leq R^2 \leq 1.$$

- Suppose you want to predict  $Y$ . If you do not know the  $x$ 's, then the best prediction is  $\bar{y}$ . The variability corresponding to this prediction is expressed by the *total variation*.
- If the model is utilized for the prediction, then the prediction error is reduced to the *residual variation*.
- $R^2$  expresses the fraction of the total variation that is explained by the model.
- As more variables are added to the model,  $D(\mathbf{y}; p_1(\mathbf{y}))$  will decrease, and  $R^2$  will increase.

Adjusted coefficient of determination,  $R^2_{adj}$ 

- The *adjusted coefficient of determination* aims to correct that  $R^2$  increases as more variables are added to the model.
- It is defined as:

$$R^2_{adj} = 1 - \frac{D(\mathbf{y}; p_1(\mathbf{y})) / (n - k)}{D(\mathbf{y}; p_{null}(\mathbf{y})) / (n - 1)}.$$

- It charges a penalty for the number of variables in the model.
- As more variables are added to the model,  $D(\mathbf{y}; p_1(\mathbf{y}))$  decreases, but the corresponding degrees of freedom also decreases.
- The numerator in may increase if the reduction in the residual deviance caused by the additional variables does not compensate for the loss in the degrees of freedom.

# Overview

- 1 The general linear model - intro
- 2 The multivariate normal distribution and norms
- 3 Likelihood, score function and information matrix
- 4 Estimation, fitted values and residuals
- 5 Partitioning of variation and hypothesis tests
- 6 Coefficient of determination,  $R^2$
- 7 Appendix (A proof of th. 3.5)

# Proof (and more) of th. 3.5

We consider the model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}; \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \boldsymbol{\Sigma}) \quad (1)$$

The estimate of  $\boldsymbol{\beta}$  is

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{Y} \quad (2)$$

hence the expected value of  $\mathbf{Y}$  can be written as

$$\hat{\mathbf{Y}} = \mathbf{X}(\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{Y} = \mathbf{H}_y \mathbf{Y} \quad (3)$$

If  $\boldsymbol{\Sigma} = \mathbf{I}$  then  $\mathbf{H}_y$  is a projection matrix and the proof can be based on properties of projection matrices. However in the general case  $\mathbf{H}_y$  is not a projection matrix.

Recal that a matrix  $\mathbf{H}$  is a projection matrix iff

$$\mathbf{H} = \mathbf{H}^T; \quad \mathbf{H}^2 = \mathbf{H} \quad (4)$$

Now

$$\mathbf{H}_y^T = (\mathbf{X}(\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Sigma}^{-1})^T \quad (5)$$

$$= ((\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Sigma}^{-1})^T \mathbf{X}^T \quad (6)$$

$$= (\mathbf{X}^T \mathbf{\Sigma}^{-1})^T (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \quad (7)$$

$$= \mathbf{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \neq \mathbf{H}_y \quad (8)$$

Hence  $\mathbf{H}_y$  is not a projection matrix.

Now consider the (non unique) “square root” ( $T$ ) of  $\Sigma = TT^T$  (since  $\Sigma$  is a variance covariance matrix this square root exist).

$$Y = X\beta + Te; \quad e \sim N(0, \sigma^2 I) \quad (9)$$

Since  $V(Te) = \sigma^2 \Sigma$  we have  $Te \sim N(0, \sigma^2 \Sigma)$ .

Now we can rewrite this as ( $Z = T^{-1}Y$ )

$$Z = T^{-1}X\beta + e; \quad e \sim N(0, \sigma^2 I) \quad (10)$$

The resulting estimate  $\tilde{\beta}$  is (we use the short hand notation  $T^{-T} = (T^{-1})^T$ )

$$\tilde{\beta} = (X^T T^{-T} T^{-1} X)^{-1} X^T T^{-T} Z \quad (11)$$

$$= (X^T \Sigma^{-1} X)^{-1} X^T T^{-T} T^{-1} Y \quad (12)$$

$$= (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y \quad (13)$$

$$= \hat{\beta} \quad (14)$$

The expected value of  $\mathbf{Z}$  is

$$\hat{\mathbf{Z}} = \mathbf{T}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T} \mathbf{Z} \quad (15)$$

$$= \mathbf{H}_z \mathbf{Z} \quad (16)$$

The matrix  $\mathbf{H}_z$  is a projection matrix

$$\mathbf{H}_z^T = (\mathbf{T}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T})^T \quad (17)$$

$$= (\mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T})^T \mathbf{T}^{-T} \quad (18)$$

$$= \mathbf{T}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T} = \mathbf{H}_z \quad (19)$$

$$\mathbf{H}_z^2 = \mathbf{T}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T} \mathbf{T}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T} \quad (20)$$

$$= \mathbf{T}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T} \quad (21)$$

$$= \mathbf{T}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T} = \mathbf{H}_z \quad (22)$$

Hence  $\mathbf{H}_z$  is a projection matrix.

Note that since  $\mathbf{H}_z$  is a projection matrix then  $\mathbf{I} - \mathbf{H}_z$  is also a projection matrix. The expected value of  $SSE_z$  is

$$E[SSE_z] = E[\mathbf{Z}^T(\mathbf{I} - \mathbf{H}_z)\mathbf{Z}] \quad (23)$$

Using prop. 3.22 of Bingham and Fry (2010), we can rewrite

$$E[SSE_z] = \text{Tr}((\mathbf{I} - \mathbf{H}_z)V[\mathbf{Z}]) + E[\mathbf{Z}^T](\mathbf{I} - \mathbf{H}_z)E[\mathbf{Z}] \quad (24)$$

$$= \sigma^2 \text{Tr}(\mathbf{I} - \mathbf{H}_z) + E[\mathbf{Z}^T](\mathbf{I} - \mathbf{H}_z)E[\mathbf{Z}] \quad (25)$$

Take it piece by piece

$$E[\mathbf{Z}^T](\mathbf{I} - \mathbf{H}_z) = \beta^T \mathbf{X}^T \mathbf{T}^{-T} (\mathbf{I} - \mathbf{T}^{-1} \mathbf{X} (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T}) \quad (26)$$

$$= \beta^T (\mathbf{X}^T \mathbf{T}^{-T} - \mathbf{X}^T \Sigma^{-1} \mathbf{X} (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T}) \quad (27)$$

$$= \beta^T \mathbf{0} = \mathbf{0} \quad (28)$$

Hence

$$E[SSE_z] = \sigma^2 (n - \text{Tr}(\mathbf{H}_z)) \quad (29)$$



The second piece

$$Tr(\mathbf{H}_z) = Tr(\mathbf{T}^{-1} \mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T}) \quad (30)$$

$$= Tr(\mathbf{X} (\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}^{-T} \mathbf{T}^{-1}) \quad (31)$$

$$= Tr((\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X}) \quad (32)$$

$$= Tr(\mathbf{I}_p) = p \quad (33)$$

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

$$E[SSE_z] = \sigma^2(n - p) \quad (34)$$

Also

$$SSE_z = (\mathbf{T}^{-1} \mathbf{Y} - \mathbf{T}^{-1} \mathbf{X} \hat{\beta})^T (\mathbf{T}^{-1} \mathbf{Y} - \mathbf{T}^{-1} \mathbf{X} \hat{\beta}) \quad (35)$$

$$= (\mathbf{Y} - \mathbf{X} \hat{\beta})^T \mathbf{T}^{-T} \mathbf{T}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\beta}) \quad (36)$$

$$= (\mathbf{Y} - \mathbf{X} \hat{\beta})^T \mathbf{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X} \hat{\beta}) \quad (37)$$

$$= SSE_y \quad (38)$$

Hence an unbiased estimator for  $\sigma^2$  is

$$\hat{\sigma}^2 = \frac{SSE_z}{n - p} \quad (39)$$

$$= \frac{SSE_y}{n - p} \quad (40)$$

$$= \frac{(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{n - p} \quad (41)$$

Note also that the distribution  $(\sigma^2 \chi_{n-p}^2)$  of  $\hat{\sigma}^2$  follows directly from  $SSE_z$