

Introduction to General and Generalized Linear Models

General Linear Models - part I

Henrik Madsen
Poul Thyregod

Informatics and Mathematical Modelling
Technical University of Denmark
DK-2800 Kgs. Lyngby

February 2012

Today

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

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.

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, σ^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

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

Deviance for normal distributed variables

Definition (Deviance for normal distributed variables)

Let us introduce the notation

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

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, δ_{Σ} :

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

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):

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

The score function, observed - and expected information for μ

- The score function wrt. μ is

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

- The observed information (wrt. μ) is

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

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

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

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

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

General Linear Models

- In the *linear model* it is assumed that $\boldsymbol{\mu}$ belongs to a linear (or affine) subspace Ω_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. $\hat{\boldsymbol{\mu}} = \mathbf{y}$.
- 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 μ . The hypothesis is of the form

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

where Ω_0 is a linear subspace of \mathbb{R}^n of dimension k , and where μ_0 denotes a vector of *known offset values*.

Definition (Dimension of general linear model)

The dimension of the subspace Ω_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

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

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.

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 δ_{Σ}), $p_0(\mathbf{y})$ of \mathbf{y} onto the linear subspace Ω_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 Σ

Notice that it has been assumed that Σ is known. If Σ is unknown, one possibility is to use the relaxation algorithm described in Madsen (2008)^a.

^aMadsen, H. (2008) Time Series Analysis. Chapman, Hall

Fitted values

Fitted – or predicted – values

The *fitted* values $\hat{\boldsymbol{\mu}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ is found as the projection of \mathbf{y} (denoted $p_0(\mathbf{y})$) on to the subspace Ω_0 spanned by \mathbf{X} , and $\hat{\boldsymbol{\beta}}$ 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

- The matrix

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

is a projection matrix.

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

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

- It follows that the predicted values are normally distributed with

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

- 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 μ 's.

Residuals

The observed *residuals* are

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

Orthogonality

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

The normal equations show that

$$\mathbf{X}^T \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0}$$

i.e. the *residuals* are orthogonal (with respect to $\boldsymbol{\Sigma}^{-1}$) to the subspace Ω_0 .

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

Residuals

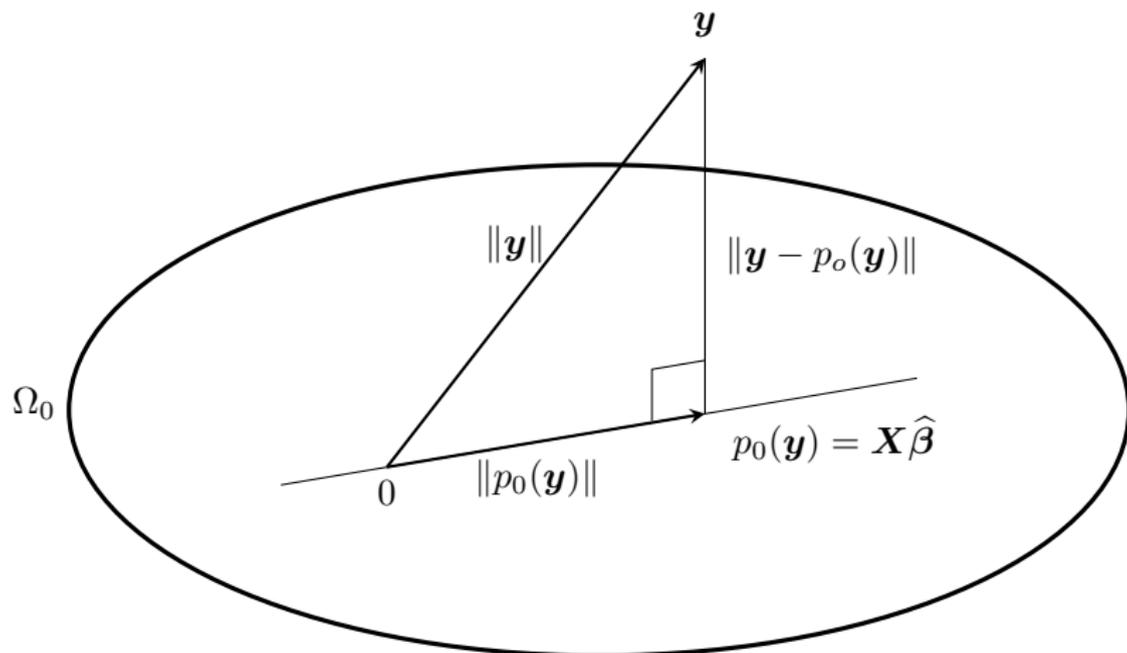


Figure: Orthogonality between the residual $(\mathbf{y} - \mathbf{X}\hat{\beta})$ and the vector $\mathbf{X}\hat{\beta}$.

Residuals

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

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

- 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 Ω_0 .
- It may be shown that the distribution of the residuals \mathbf{r} is independent of the fitted values $\mathbf{X}\hat{\boldsymbol{\beta}}$.

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)

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

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

χ^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 χ^2 distributed variables with $n - k$ and k degrees of freedom, respectively.

Estimation of the residual variance σ^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 σ^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 $f = n - k$.

Summary: General Linear Model

- A general linear model is:

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

Consider the well known two way ANOVA:

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

An expanded view of this model is:

$$\begin{array}{rcllcl}
 y_{11} & = & \mu & + & \alpha_1 & & + & \beta_1 & & + & \varepsilon_{11} \\
 y_{21} & = & \mu & & & + & \alpha_2 & + & \beta_1 & & + & \varepsilon_{21} \\
 y_{12} & = & \mu & + & \alpha_1 & & & & + & \beta_2 & & + & \varepsilon_{12} \\
 y_{22} & = & \mu & & & + & \alpha_2 & & & + & \beta_2 & & + & \varepsilon_{22} \\
 y_{13} & = & \mu & + & \alpha_1 & & & & & + & \beta_3 & & + & \varepsilon_{13} \\
 y_{23} & = & \mu & & & + & \alpha_2 & & & + & \beta_3 & & + & \varepsilon_{23}
 \end{array} \tag{1}$$

The exact same in matrix notation:

$$\underbrace{\begin{pmatrix} y_{11} \\ y_{21} \\ y_{12} \\ y_{22} \\ y_{13} \\ y_{23} \end{pmatrix}}_{\mathbf{y}} = \underbrace{\begin{pmatrix} 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 & 1 \end{pmatrix}}_{\mathbf{X}} \underbrace{\begin{pmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}}_{\boldsymbol{\beta}} + \underbrace{\begin{pmatrix} \varepsilon_{11} \\ \varepsilon_{21} \\ \varepsilon_{12} \\ \varepsilon_{22} \\ \varepsilon_{13} \\ \varepsilon_{23} \end{pmatrix}}_{\boldsymbol{\varepsilon}} \tag{2}$$

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

- \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} .

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:

$$\mathbf{y} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ 1 & x_n \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + \varepsilon$$

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 Ω_1 , and hence we introduce $\Omega_0 \subset \Omega_1$ as a linear (affine) subspace with $\dim(\Omega_0) = m_0$.

Model reduction

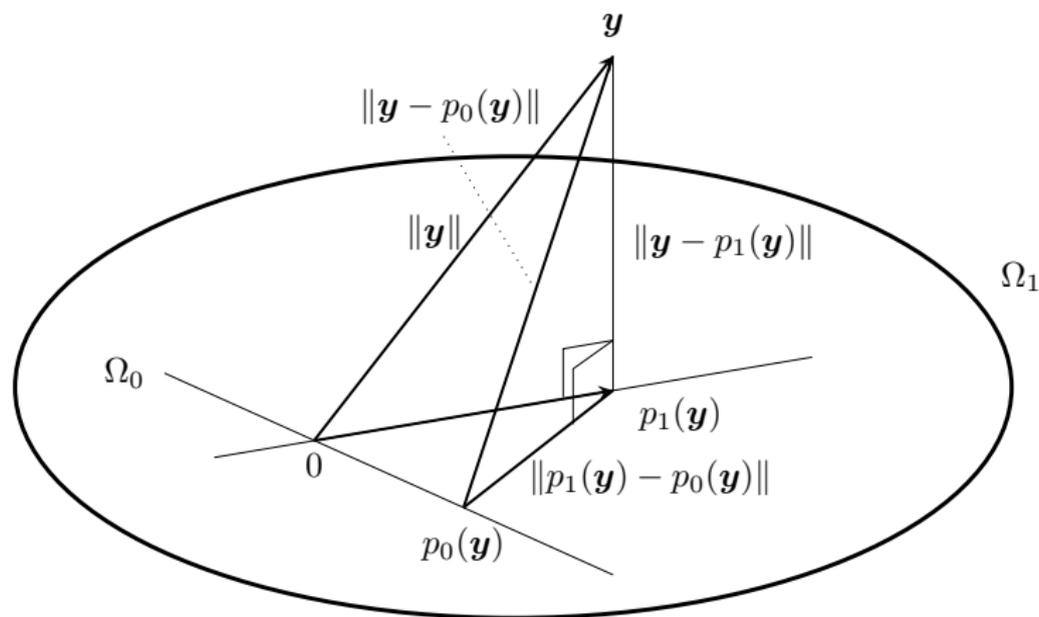


Figure: Model reduction. The partitioning of the deviance corresponding to a test of the hypothesis $\mathcal{H}_0 : \mu \in \Omega_0$ under the assumption of $\mathcal{H}_1 : \mu \in \Omega_1$.

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 Ω_1 and Ω_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 χ^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 σ^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 Ω_1 and Ω_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 $p = P[F(k - 1, n - k) \geq F_{obs}]$.

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_{adj}^2

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

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

- 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.