## 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

October 2010

### 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[Y] = 0 and the variance-covariance matrix Var[Y] = I.

#### Definition (Multivariate normal distribution)

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

Since 
$$A$$
 and  $b$  are fixed, we have  $E[Z] = b$  and  $Var[Z] = AA^T$ .

### The multivariate normal distribution

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

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

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

where  $\Sigma$  is seen to be (a) symmetric and (b) positive semi-definite. We write  $Z \sim N_k(\mu, \sigma^2 \Sigma)$ .

### 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

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

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

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

### Inner product and norm

#### Definition (Inner product and norm)

The bilinear form

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

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

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

### Deviance for normal distributed variables

Definition (Deviance for normal distributed variables)

Let us introduce the notation

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

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

For a normal distribution with  $\Sigma = I$ , 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,  $\delta_{\Sigma}$ :

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

### The likelihood and log-likelihood function

• The likelihood function is:

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

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

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

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

• The score function wrt.  $\mu$  is

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

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

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

• It is seen that the observed information does not depend on the observations y. Hence the expected information is

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

### 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  $(Y_1, Y_2, \ldots, Y_n)$  can be written on the matrix form

$$oldsymbol{Y} = oldsymbol{\mu} + oldsymbol{\epsilon}$$

where

$$\boldsymbol{Y} \sim \mathrm{N}_n(\boldsymbol{\mu}, \sigma^2 \boldsymbol{\Sigma}) \ ext{for} \ \boldsymbol{y} \in \mathbb{R}^n$$

### General Linear Models

- In the *linear model* it is assumed that μ belongs to a linear (or affine) subspace Ω<sub>0</sub> of ℝ<sup>n</sup>.
- 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  $Y_1, Y_2, \ldots, Y_n$  is normally distributed as described before. A general linear model for  $Y_1, Y_2, \ldots, Y_n$  is a model where an affine hypothesis is formulated for  $\mu$ . The hypothesis is of the form

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

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

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

$$\mathcal{H}_0: oldsymbol{\mu} - oldsymbol{\mu}_0 = oldsymbol{X}oldsymbol{eta}$$
 with  $oldsymbol{eta} \in \mathbb{R}^k,$ 

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

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

$$oldsymbol{x}_i^T = \left(egin{array}{c} x_{i1}\ x_{i2}\ dots\ x_{ik}\ \end{array}
ight)^T,$$

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

Theorem (ML estimates of mean value parameters)

For hypothesis of the form

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

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

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

If X has full rank, the solution is uniquely given by

$$\widehat{oldsymbol{eta}} = (oldsymbol{X}^T oldsymbol{\Sigma}^{-1} oldsymbol{X})^{-1} oldsymbol{X}^T oldsymbol{\Sigma}^{-1} oldsymbol{y}$$

### Properties of the ML estimator

Theorem (Properties of the ML estimator)

For the ML estimator we have

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

#### Unknown $\Sigma$

Notice that it has been assumed that  $\Sigma$  is known. If  $\Sigma$  is unknown, one possibility is to use the relaxation algorithm described in Madsen (2008)<sup>*a*</sup>.

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

### Fitted values

#### Fitted – or predicted – values

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

#### Definition (Projection matrix)

A matrix H is a *projection matrix* if and only if (a)  $H^T = H$  and (b)  $H^2 = H$ , i.e. the matrix is *idempotent*.

### The hat matrix

The matrix

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

is a projection matrix.

• The projection matrix provides the predicted values  $\widehat{\mu}$ , since

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

• It follows that the predicted values are normally distributed with

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

The matrix *H* is often termed the *hat matrix* since it transforms the observations *y* to their predicted values symbolized by a "hat" on the μ's.

### Residuals

The observed residuals are

$$r = y - X\widehat{eta} = (I - H)y$$

#### Orthogonality

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

The normal equations show that

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

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

### Residuals



Figure: Orthogonality between the residual  $(\boldsymbol{y}-\boldsymbol{X}\widehat{eta})$  and the vector  $\boldsymbol{X}\widehat{eta}.$ 

### Residuals

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

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

- 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 r is independent of the fitted values  $X\widehat{eta}$ .

### 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)

$$Y^T Y = Y^T H_1 Y + Y^T H_2 Y + \dots + Y^T H_k Y$$

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

- i The ranks of the  $oldsymbol{H}_i$  adds to n, i.e.  $\sum_{i=1}^k n_i = n$
- ii Each quadratic form  $\mathbf{Y}^T \mathbf{H}_i \mathbf{Y} \sim \chi^2_{n_i}$  (thus the  $H_i$  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

$$\begin{split} \mathrm{D}(\boldsymbol{y};\boldsymbol{X}\boldsymbol{\beta}) &= \mathrm{D}(\boldsymbol{y};\boldsymbol{X}\widehat{\boldsymbol{\beta}}) + \mathrm{D}(\boldsymbol{X}\widehat{\boldsymbol{\beta}};\boldsymbol{X}\boldsymbol{\beta}) \\ &= (\boldsymbol{y}-\boldsymbol{X}\widehat{\boldsymbol{\beta}})^T\boldsymbol{\Sigma}^{-1}(\boldsymbol{y}-\boldsymbol{X}\widehat{\boldsymbol{\beta}}) \\ &+ (\widehat{\boldsymbol{\beta}}-\boldsymbol{\beta})^T\boldsymbol{X}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{X}(\widehat{\boldsymbol{\beta}}-\boldsymbol{\beta}) \\ &\geq (\boldsymbol{y}-\boldsymbol{X}\widehat{\boldsymbol{\beta}})^T\boldsymbol{\Sigma}^{-1}(\boldsymbol{y}-\boldsymbol{X}\widehat{\boldsymbol{\beta}}) \end{split}$$

### Partitioning of variation

### $\chi^2\text{-distribution}$ of individual contributions

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

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

Furthermore, it follows from the normal distribution of r and of  $\widehat{oldsymbol{eta}}$  that

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

moreover, the independence of r and  $X\widehat{\beta}$  implies that  $D(y; X\widehat{\beta})$  and  $D(X\widehat{\beta}; X\beta)$  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

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

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

$$\widehat{\sigma}^2 = rac{\mathrm{D}(oldsymbol{y};oldsymbol{X}\widehat{oldsymbol{eta}})}{n-k} = rac{(oldsymbol{y}-oldsymbol{X}\widehat{oldsymbol{eta}})^T \mathbf{\Sigma}^{-1}(oldsymbol{y}-oldsymbol{X}\widehat{oldsymbol{eta}})}{n-k}$$

Under the hypothesis,  $\widehat{\sigma}^2 \sim \sigma^2 \chi_f^2 / f$  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  $\boldsymbol{Y} \sim \mathrm{N}_n(\boldsymbol{\mu},\sigma^2\boldsymbol{\Sigma}).$
- 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  $\mu$  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

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

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  $p_1(\boldsymbol{y})$  and  $p_0(\boldsymbol{y})$  denote the projection of  $\boldsymbol{y}$  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  $p = P[F(m_1 - m_0, n - m_1) \ge F_{obs}]$ , 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(m{y}) - p_0(m{y})  ^2$	$\frac{  p_1(\boldsymbol{y}) - p_0(\boldsymbol{y})  ^2/(m_1 - m_0)}{  \boldsymbol{y} - p_1(\boldsymbol{y})  ^2/(n - m_1)}$
Residual under model	$n-m_1$	$  m{y} - p_1(m{y})  ^2$	
Residual under hypothesis	$n-m_0$	$   m{y} - p_0(m{y})  ^2$	

Table: Deviance table corresponding to a test for model reduction as specified by  $\mathcal{H}_0$ . For  $\Sigma = I$  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  $\mathcal{H}_0: \mu \in \Omega_0$  under the assumption that  $\mathcal{H}_1: \mu \in \Omega_1$  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  $D(\boldsymbol{y}; p_1(\boldsymbol{y}))$ .

# 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 Ω<sub>null</sub> = ℝ, i.e. dim Ω<sub>null</sub> = 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(\boldsymbol{y}) - p_{null}(\boldsymbol{y})  ^2$	$\frac{  p_1(\boldsymbol{y}) - p_{null}(\boldsymbol{y})  ^2/(k-1)}{  u_1 - u_2(\boldsymbol{x})  ^2/(k-1)}$
Residual under $\mathcal{H}_1$	n-k	$  oldsymbol{y}-p_1(oldsymbol{y})  ^2$	$  y - p_1(y)  ^2 / (n - \kappa)$
Total	n-1	$  oldsymbol{y}-p_{null}(oldsymbol{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

$$R^2 = \frac{\mathrm{D}(p_1(\boldsymbol{y}); p_{null}(\boldsymbol{y}))}{\mathrm{D}(\boldsymbol{y}; p_{null}(\boldsymbol{y}))} = 1 - \frac{\mathrm{D}(\boldsymbol{y}; p_1(\boldsymbol{y}))}{\mathrm{D}(\boldsymbol{y}; p_{null}(\boldsymbol{y}))}, \ \ 0 \le R^2 \le 1.$$

- Suppose you want to predict Y. If you do not know the x's, then the best prediction is  $\overline{y}$ . 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,  $D(y; p_1(y))$  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:

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

- It charges a penalty for the number of variables in the model.
- As more variables are added to the model,  $D(y; p_1(y))$  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.