

Introduction to General and Generalized Linear Models

Generalized Linear Models - part II

Henrik Madsen
Poul Thyregod

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

October 2010

Today

- The generalized linear model
 - Link function
 - (Estimation)
 - Fitted values
 - Residuals
- Likelihood ratio test
- Over-dispersion

The Generalized Linear Model

Definition (The generalized linear model)

Assume that Y_1, Y_2, \dots, Y_n are mutually independent, and the density can be described by an exponential dispersion model with the same variance function $V(\mu)$.

A *generalized linear model* for Y_1, Y_2, \dots, Y_n describes an affine hypothesis for $\eta_1, \eta_2, \dots, \eta_n$, where

$$\eta_i = g(\mu_i)$$

is a transformation of the mean values $\mu_1, \mu_2, \dots, \mu_n$.

The hypothesis is of the form

$$\mathcal{H}_0 : \boldsymbol{\eta} - \boldsymbol{\eta}_0 \in L,$$

where L is a linear subspace \mathbb{R}^n of dimension k , and where $\boldsymbol{\eta}_0$ denotes a vector of *known off-set values*.

Dimension and design matrix

Definition (Dimension of the generalized linear model)

The dimension k of the subspace L for the generalized linear model is the *dimension of the model*

Definition (Design matrix for the generalized linear model)

Consider the linear subspace $L = \text{span}\{x_1, \dots, x_k\}$, i.e. the subspace is spanned by k vectors ($k < n$), such that the hypothesis can be written

$$\eta - \eta_0 = \mathbf{X}\boldsymbol{\beta} \text{ with } \boldsymbol{\beta} \in \mathbb{R}^k,$$

where \mathbf{X} has full rank. The $n \times k$ matrix \mathbf{X} is called the *design matrix*. The i^{th} row of the design matrix is given by the *model vector*

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

for the i^{th} observation.

The link function

Definition (The link function)

The *link function*, $g(\cdot)$ describes the relation between the linear predictor η_i and the mean value parameter $\mu_i = E[Y_i]$. The relation is

$$\eta_i = g(\mu_i)$$

The inverse mapping $g^{-1}(\cdot)$ thus expresses the mean value μ as a function of the linear predictor η :

$$\mu = g^{-1}(\eta)$$

that is

$$\mu_i = g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}) = g^{-1} \left(\sum_j x_{ij} \beta_j \right)$$

Link functions

The most commonly used link functions, $\eta = g(\mu)$, are :

Name	Link function $\eta = g(\mu)$	$\mu = g^{-1}(\eta)$
Identity	μ	η
logarithm	$\ln(\mu)$	$\exp(\eta)$
logit	$\ln(\mu/(1 - \mu))$	$\exp(\eta)/[1 + \exp(\eta)]$
reciprocal	$1/\mu$	$1/\eta$
power	μ^k	$\eta^{1/k}$
squareroot	$\sqrt{\mu}$	η^2
probit	$\Phi^{-1}(\mu)$	$\Phi(\eta)$
log-log	$\ln(-\ln(\mu))$	$\exp(-\exp(\eta))$
cloglog	$\ln(-\ln(1 - \mu))$	$1 - \exp(-\exp(\eta))$

Table: Commonly used link function.

The canonical link

The canonical link is the function which transforms the mean to the canonical location parameter of the exponential dispersion family, i.e. it is the function for which $g(\mu) = \theta$. The canonical link function for the most widely considered densities are

Density	Link: $\eta = g(\mu)$	Name
Normal	$\eta = \mu$	identity
Poisson	$\eta = \ln(\mu)$	logarithm
Binomial	$\eta = \ln[\mu/(1 - \mu)]$	logit
Gamma	$\eta = 1/\mu$	reciprocal
Inverse Gauss	$\eta = 1/\mu^2$	power ($k = -2$)

Table: Canonical link functions for some widely used densities.

Specification of a generalized linear model

- a) Distribution / Variance function:

Specification of the distribution – or the *variance function* $V(\mu)$.

- b) Link function:

Specification of the *link function* $g(\cdot)$, which describes a function of the mean value which can be described linearly by the explanatory variables.

- c) Linear predictor:

Specification of the linear dependency

$$g(\mu_i) = \eta_i = (\mathbf{x}_i)^T \boldsymbol{\beta}.$$

- d) Precision (optional):

If needed the precision is formulated as *known individual weights*, $\lambda_i = w_i$, or as a *common dispersion parameter*, $\lambda = 1/\sigma^2$, or a *combination* $\lambda_i = w_i/\sigma^2$.

Maximum likelihood estimation

Theorem (Estimation in generalized linear models)

Consider the generalized linear model as defined on slide 3 for the observations Y_1, \dots, Y_n and assume that Y_1, \dots, Y_n are mutually independent with densities, which can be described by an exponential dispersion model with the variance function $V(\cdot)$, dispersion parameter σ^2 , and optionally the weights w_i .

Assume that the linear predictor is parameterized with β corresponding to the design matrix \mathbf{X} , then the maximum likelihood estimate $\hat{\beta}$ for β is found as the solution to

$$[\mathbf{X}(\beta)]^T \mathbf{i}_\mu(\mu)(\mathbf{y} - \mu) = \mathbf{0},$$

where $\mathbf{X}(\beta)$ denotes the local design matrix and $\mu = \mu(\beta)$ given by

$$\mu_i(\beta) = g^{-1}(\mathbf{x}_i^T \beta),$$

denotes the fitted mean values corresponding to the parameters β , and $\mathbf{i}_\mu(\mu)$ is the expected information with respect to μ .

Properties of the ML estimator

Theorem (Asymptotic distribution of the ML estimator)

Under the hypothesis $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$ we have asymptotically

$$\frac{\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}}{\sqrt{\sigma^2}} \in N_k(\mathbf{0}, \boldsymbol{\Sigma}),$$

where the dispersion matrix $\boldsymbol{\Sigma}$ for $\hat{\boldsymbol{\beta}}$ is

$$D[\hat{\boldsymbol{\beta}}] = \boldsymbol{\Sigma} = [\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}) \mathbf{X}]^{-1}$$

with

$$\mathbf{W}(\boldsymbol{\beta}) = \text{diag} \left\{ \frac{w_i}{[g'(\mu_i)]^2 V(\mu_i)} \right\},$$

In the case of the canonical link, the weight matrix $\mathbf{W}(\boldsymbol{\beta})$ is

$$\mathbf{W}(\boldsymbol{\beta}) = \text{diag} \{w_i V(\mu_i)\}.$$

Linear prediction for the generalized linear model

Definition (Linear prediction for the generalized linear model)

The linear prediction $\hat{\boldsymbol{\eta}}$ is defined as the values

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

with the linear prediction corresponding to the i 'th observation is

$$\hat{\eta}_i = \sum_{j=1}^k x_{ij}\hat{\beta}_j = (\mathbf{x}_i)^T \hat{\boldsymbol{\beta}}.$$

The linear predictions $\hat{\boldsymbol{\eta}}$ are approximately normally distributed with

$$D[\hat{\boldsymbol{\eta}}] \approx \hat{\sigma}^2 \mathbf{X}\boldsymbol{\Sigma}\mathbf{X}^T$$

where $\boldsymbol{\Sigma}$ is the dispersion matrix for $\hat{\boldsymbol{\beta}}$.

Fitted values for the generalized linear model

Definition (Fitted values for the generalized linear model)

The fitted values are defined as the values

$$\hat{\boldsymbol{\mu}} = \boldsymbol{\mu}(\mathbf{X}\hat{\boldsymbol{\beta}}),$$

where the i^{th} value is given as

$$\hat{\mu}_i = g^{-1}(\hat{\eta}_i)$$

with the fitted value $\hat{\eta}_i$ of the linear prediction.

The fitted values $\hat{\boldsymbol{\mu}}$ are approximately normally distributed with

$$D[\hat{\boldsymbol{\mu}}] \approx \hat{\sigma}^2 \left[\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} \right]^2 \mathbf{X}\boldsymbol{\Sigma}\mathbf{X}^T$$

where $\boldsymbol{\Sigma}$ is the dispersion matrix for $\hat{\boldsymbol{\beta}}$.

Residual deviance

Definition (Residual deviance)

Consider the generalized linear model defined on slide 3. The *residual deviance* corresponding to this model is

$$D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}})) = \sum_{i=1}^n w_i d(y_i; \hat{\mu}_i)$$

with $d(y_i; \hat{\mu}_i)$ denoting the unit deviance corresponding the observation y_i and the fitted value $\hat{\mu}_i$ and where w_i denotes the weights (if present).

If the model includes a dispersion parameter σ^2 , the *scaled* residual deviance is

$$D^*(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}})) = \frac{D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))}{\sigma^2}.$$

Residuals

Residuals represents the difference between the data and the model. In the classical GLM the residuals are $r_i = y_i - \hat{\mu}_i$. These are called response residuals for GLM's. Since the variance of the response is not constant for most GLM's we need some modification. We will look at:

- Deviance residuals
- Pearson residuals

Residuals

Definition (Deviance residual)

Consider the generalized linear model from for the observations Y_1, \dots, Y_n .

The *deviance residual* for the i 'th observation is defined as

$$r_i^D = r_D(y_i; \hat{\mu}_i) = \text{sign}(y_i - \hat{\mu}_i) \sqrt{w_i d(y_i, \hat{\mu}_i)}$$

where $\text{sign}(x)$ denotes the *sign function* $\text{sign}(x) = 1$ for $x > 0$ og $\text{sign}(x) = -1$ for $x < 0$, and with w_i denoting the weight (if relevant), $d(y; \mu)$ denoting the unit deviance and $\hat{\mu}_i$ denoting the fitted value corresponding to the i 'th observation.

Assessments of the deviance residuals is in good agreement with the likelihood approach as the deviance residuals simply express differences in log-likelihood.

Residuals

Definition (Pearson residual)

Consider again the generalized linear model from for the observations Y_1, \dots, Y_n .

The *Pearson residuals* are defined as the values

$$r_i^P = r_P(y_i; \hat{\mu}_i) = \frac{y_i - \hat{\mu}_i}{\sqrt{V(\hat{\mu}_i)/w_i}}$$

The Pearson residual is thus obtained by scaling the response residual with $\sqrt{\text{Var}[Y_i]}$. Hence, the Pearson residual is the response residual normalized with the estimated standard deviation for the observation.

Likelihood ratio tests

- The approximative normal distribution of the ML-estimator implies that many distributional results from the classical GLM-theory are carried over to generalized linear models as approximative (asymptotic) results.
- An example of this is the likelihood ratio test.
- In the classical GLM case it was possible to derive the exact distribution of the likelihood ratio test statistic (the F-distribution).
- For generalized linear models, this is not possible, and hence we shall use the asymptotic results for the logarithm of the likelihood ratio.

Likelihood ratio test

Theorem (Likelihood ratio test)

Consider the generalized linear model. Assume that the model

$$\mathcal{H}_1 : \boldsymbol{\eta} \in L \subset \mathbb{R}^k$$

holds with L parameterized as $\boldsymbol{\eta} = \mathbf{X}_1\boldsymbol{\beta}$, and consider the hypotheses

$$\mathcal{H}_0 : \boldsymbol{\eta} \in L_0 \subset \mathbb{R}^m$$

where $\boldsymbol{\eta} = \mathbf{X}_0\boldsymbol{\alpha}$ and $m < k$, and with the alternative $\mathcal{H}_1 : \boldsymbol{\eta} \in L \setminus L_0$. Then the likelihood ratio test for \mathcal{H}_0 has the test statistic

$$-2 \log \lambda(\mathbf{y}) = D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}})) - D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\alpha}}))$$

When \mathcal{H}_0 is true, the test statistic will asymptotically follow a $\chi^2(k - m)$ distribution.

If the model includes a dispersion parameter, σ^2 , then $D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}); \boldsymbol{\mu}(\boldsymbol{\beta}(\hat{\boldsymbol{\alpha}})))$ will asymptotically follow a $\sigma^2\chi^2(k - m)$ distribution.

Test for model 'sufficiency'

- In analogy with classical GLM's one often starts with formulating a rather comprehensive model, and then reduces the model by successive tests.
- In contrast to classical GLM's we may however test the goodness of fit of the initial model.
- The test is a special case of the likelihood ratio test.

Test for model 'sufficiency'

Test for model 'sufficiency'

Consider the generalized linear model, and assume that the dispersion $\sigma^2 = 1$.

Let \mathcal{H}_{full} denote the *full*, or *saturated* model, i.e. $\mathcal{H}_{full} : \boldsymbol{\mu} \in \mathbb{R}^n$ and consider the hypotheses

$$\mathcal{H}_0 : \boldsymbol{\eta} \in L \subset \mathbb{R}^k$$

with L parameterized as $\boldsymbol{\eta} = \mathbf{X}_0\boldsymbol{\beta}$.

Then, as the residual deviance under \mathcal{H}_{full} is 0, the test statistic is the residual deviance $D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$. When \mathcal{H}_0 is true, the test statistic is distributed as $\chi^2(n - k)$. The test rejects for large values of $D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$.

Residual deviance measures goodness of fit

- The residual deviance $D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$ is a reasonable measure of the goodness of fit of a model \mathcal{H}_0 .
- When referring to a hypothesized model \mathcal{H}_0 , we shall sometimes use the symbol $G^2(\mathcal{H}_0)$ to denote the residual deviance $D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$.
- Using that convention, the partitioning of residual deviance may be formulated as

$$G^2(\mathcal{H}_0|\mathcal{H}_1) = G^2(\mathcal{H}_0) - G^2(\mathcal{H}_1)$$

with $G^2(\mathcal{H}_0|\mathcal{H}_1)$ interpreted as the goodness fit test statistic for \mathcal{H}_0 conditioned on \mathcal{H}_1 being true, and $G^2(\mathcal{H}_0)$ and $G^2(\mathcal{H}_1)$, denoting the unconditional goodness of fit statistics for \mathcal{H}_0 and \mathcal{H}_1 , respectively.

Analysis of deviance table

- The initial test for goodness of fit of the initial model is often represented in an *analysis of deviance table* in analogy with the ANOVA table for classical GLM's.
- In the table the goodness of fit test statistic corresponding to the initial model $G^2(\mathcal{H}_1) = D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$ is shown in the line labelled "Error".
- The statistic should be compared to percentiles in the $\chi^2(n - k)$ distribution.
- The table also shows the test statistic for \mathcal{H}_{null} under the assumption that \mathcal{H}_1 is true.
- The test investigates whether the model is necessary at all, i.e. whether at least some of the coefficients differ significantly from zero.

Analysis of deviance table

- Note, that in the case of a generalized linear model, we can start the analysis by using the residual (error) deviance to test whether the model may be maintained, at all.
- This is in contrast to the classical GLM's where the residual sum of squares around the initial model \mathcal{H}_1 served to estimate σ^2 , and therefore we had no reference value to compare with the residual sum of squares.
- In the generalized linear models the variance is a known function of the mean, and therefore in general there is no need to estimate a separate variance.

Analysis of deviance table

Source	f	Deviance	Mean deviance	Goodness of fit interpretation
Model \mathcal{H}_{null}	$k - 1$	$D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}); \hat{\boldsymbol{\mu}}_{null})$	$\frac{D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}); \hat{\boldsymbol{\mu}}_{null})}{k - 1}$	$G^2(\mathcal{H}_{null} \mathcal{H}_1)$
Residual (Error)	$n - k$	$D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$	$\frac{D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))}{n - k}$	$G^2(\mathcal{H}_1)$
Corrected total	$n - 1$	$D(\mathbf{y}; \hat{\boldsymbol{\mu}}_{null})$		$G^2(\mathcal{H}_{null})$

Table: Initial assessment of goodness of fit of a model \mathcal{H}_0 . \mathcal{H}_{null} and $\hat{\boldsymbol{\mu}}_{null}$ refer to the *minimal model*, i.e. a model with all observations having the same mean value.

Overdispersion

- It may happen that even if one has tried to fit a rather comprehensive model (i.e. a model with many parameters), the fit is not satisfactory, and the residual deviance $D(\mathbf{y}; \mu(\hat{\beta}))$ is larger than what can be explained by the χ^2 -distribution.
- An explanation for such a poor model fit could be an improper choice of linear predictor, or of link or response distribution.
- If the residuals exhibit a random pattern, and there are no other indications of misfit, then the explanation could be that the variance is larger than indicated by $V(\mu)$.
- We say that the data are *overdispersed*.

Overdispersion

- When data are *overdispersed*, a more appropriate model might be obtained by including a *dispersion parameter*, σ^2 , in the model, i.e. a distribution model of the form with $\lambda_i = w_i/\sigma^2$, and σ^2 denoting the overdispersion, $\text{Var}[Y_i] = \sigma^2 V(\mu_i)/w_i$.
- As the dispersion parameter only would enter in the score function as a constant factor, this does not affect the estimation of the mean value parameters β .
- However, because of the larger error variance, the distribution of the test statistics will be influenced.
- If, for some reasons, the parameter σ^2 had been known beforehand, one would include this known value in the weights, w_i .
- Most often, when it is found necessary to choose a model with overdispersion, σ^2 shall be estimated from the data.

The dispersion parameter

- For the normal distribution family, the dispersion parameter is just the variance σ^2 .
- In the case of a gamma distribution family, the shape parameter α acts as dispersion parameter.
- The maximum likelihood estimation of the shape parameter is not too complicated for the normal and the gamma distributions but for other exponential dispersion families, ML estimation of the dispersion parameter is more tricky.
- The problem is that the dispersion parameter enters in the likelihood function, not only as a factor to the deviance, but also in the normalizing factor $a(y_i, w_i/\sigma^2)$.
- It is necessary to have an explicit expression for this factor as function of σ^2 (as in the case of the normal and the gamma distribution families) in order to perform the maximum likelihood estimation.

The dispersion parameter

Approximate moment estimate for the dispersion parameter

It is common practice to use the residual deviance $D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$ as basis for the estimation of σ^2 and use the result that $D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$ is approximately distributed as $\sigma^2 \chi^2(n - k)$. It then follows that

$$\hat{\sigma}_{dev}^2 = \frac{D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))}{n - k}$$

is asymptotically unbiased for σ^2 .

Alternatively, one would utilize the corresponding Pearson goodness of fit statistic

$$X^2 = \sum_{i=1}^n w_i \frac{(y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)}$$

which likewise follows a $\sigma^2 \chi^2(n - k)$ -distribution, and use the estimator

$$\hat{\sigma}_{Pearson}^2 = \frac{X^2}{n - k}.$$

Deviance table in the case of overdispersion

Source	f	Deviance	Scaled deviance
Model \mathcal{H}_{null}	$k - 1$	$D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}); \hat{\boldsymbol{\mu}}_{null})$	$\frac{D(\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}); \hat{\boldsymbol{\mu}}_{null}) / (k - 1)}{D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}})) / (n - k)}$
Residual (Error)	$n - k$	$D(\mathbf{y}; \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}))$	
Corrected total	$n - 1$	$D(\mathbf{y}; \hat{\boldsymbol{\mu}}_{null})$	

Table: Example of Deviance table in the case of overdispersion. It is noted that the scaled deviance is equal to the model deviance scaled by the error deviance.

- The *scaled deviance*, D^* , i.e. deviance divided by $\hat{\sigma}^2$ is used in the tests instead of the crude deviance in case of overdispersion.
- For calculation of p -values etc. the asymptotic χ^2 -distribution of the scaled deviance is used.