

TMA4315 Generalized linear models H2018

Module 5: Generalized linear models - common core

*Mette Langaas, Department of Mathematical Sciences, NTNU - with contributions from
Ingeborg Hem*

11.10.2017 [PL], 12.10.2017 [IL]

Contents

Overview	2
Learning material	2
Topics	2
GLM — three ingredients	2
Random component - exponential family	2
Systematic component - linear predictor	3
Link function - and response function	3
Canonical link	4
Likelihood inference set-up	4
Loglikelihood	4
Score function	5
Expected Fisher information matrix for the GLM and covariance for $\hat{\beta}$	5
Fisher scoring and iterated reweighted least squares (IRWLS)	5
Estimator for dispersion parameter	6
Distribution of the MLE	6
What about the distribution of $\hat{\beta}, \hat{\phi}$?	7
Hypothesis testing	7
Model assessment and model choice	7
Pearson and deviance statistic	7
AIC	7
Interactive session	8
Problem 1: Exam 2011, problem 3	8
Problem 2: December 2005, Problem 2 (modified)	8
Problem 3: Exam 2006, problem 2 (a, b, d)	8
Problem 4: Exam 2007, problem 2 a, b, c	9
Problem 5: Exam UiO December 2017, Problem 2	10
Exam questions	10
December 2015	10
Further reading	11

(Latest changes: 11.010, added links to handwritten materials and dispersion formula, 07.10.2018 first version)

Overview

Learning material

- Textbook: Fahrmeir et al (2013): Chapter 5.4, 5.8.2.
- Classnotes 27.09.2018

Additional notes (with theoretical focus):

- Exponential family from Module 1
 - Proof of E and Var for exp fam
 - Proof of two forms for F
 - Orthogonal parameters
 - IRWLS
-

Topics

- random component: exponential family
 - elements: $\theta, \phi, w, b(\theta)$
 - elements for normal, binomial, Poisson and gamma
 - properties: $E(Y) = b'(\theta)$ and $\text{Var}(Y) = b''(\theta) \frac{\phi}{w}$ (and proof)
 - systematic component= linear predictor
 - requirements: full rank of design matrix
 - link function and response function
 - link examples for normal, binomial, Poisson and gamma
 - requirements: one-to-one and twice differentiable
 - canonical link
-

- likelihood inference set-up: $\theta_i \leftrightarrow \mu_i \leftrightarrow \eta_i \leftrightarrow \beta$
- the loglikelihood
- the score function
- expected Fisher information matrix for the GLM and covariance for $\hat{\beta}$
 - what about covariance of $\hat{\beta}$ when ϕ needs to be estimated?
 - estimator for dispersion parameter
- Fisher scoring and iterated reweighted least squares (IRWLS)
- Pearson and deviance statistic
- AIC

– so, for the first time: no practical examples or data sets to be analysed!

Jump to interactive.

GLM — three ingredients

Random component - exponential family

In Module 1 we introduced distributions of the Y_i , that could be written in the form of a *univariate exponential family*

$$f(y_i | \theta_i) = \exp\left(\frac{y_i \theta_i - b(\theta_i)}{\phi} \cdot w_i + c(y_i, \phi, w_i)\right)$$

where we said that

- θ_i is called the canonical parameter and is a parameter of interest
- ϕ is called a nuisance parameter (and is not of interest to us=therefore a nuisance (plage))
- w_i is a weight function, in most cases $w_i = 1$ (NB: can not contain any unknown parameters)
- b and c are known functions.

Elements - for normal, Bernoulli, Poisson and gamma

We have seen:

Distribution	θ	$b(\theta)$	ϕ	w	$E(Y) = b'(\theta)$	$b''(\theta)$	$\text{Var}(Y) = b''(\theta)\phi/w$
normal	μ	$\frac{1}{2}\theta^2$	σ^2	1	$\mu = \theta$	1	σ^2
Bernoulli	$\ln\left(\frac{p}{1-p}\right)$	$\ln(1 + \exp(\theta))$	1	1	$p = \frac{\exp(\theta)}{1+\exp(\theta)}$	$p(1-p)$	$p(1-p)$
Poisson	$\ln \mu$	$\exp(\theta)$	1	1	$\lambda = \exp(\theta)$	λ	λ
gamma	$-\frac{1}{\mu}$	$-\ln(-\theta)$	$\frac{1}{\nu}$	1	$\mu = -1/\theta$	μ^2	μ^2/ν

Properties

$$E(Y_i) = b'(\theta_i) \text{ and } \text{Var}(Y_i) = b''(\theta_i) \frac{\phi}{w_i}$$

In class we study the handwritten proof together: Proof

$b''(\theta_i)$ is often called the variance function $v(\mu_i)$.

Systematic component - linear predictor

Nothing new - as always in this course: $\eta_i = \mathbf{x}_i^T \beta$, and we require that the $n \times p$ design matrix $\mathbf{X} = (\mathbf{x}_1^T, \mathbf{x}_2^T, \dots, \mathbf{x}_n^T)$ has full rank (which is p).

Remark: in this course we always assume that $n \gg p$.

Link function - and response function

Link function

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

Response function

$$\mu_i = h(\eta_i)$$

Examples for normal, binomial, Poisson and gamma

random component	response function and link function
normal	$h(\eta_i) = \eta_i$ and $g(\mu_i) = \mu_i$, “identity link”.
binomial	$h(\eta_i) = \frac{e^{\eta_i}}{1+e^{\eta_i}}$ and $g(\mu_i) = \ln\left(\frac{\mu_i}{1-\mu_i}\right) = \text{logit}(p_i)$. NB: $\mu_i = p_i$ in our set-up.
Poisson	$h(\eta_i) = \exp(\eta_i)$ and $g(\mu_i) = \ln(\mu_i)$, log-link.
gamma	$h(\eta_i) = -\frac{1}{\eta_i}$ and $g(\mu_i) = -\frac{1}{\mu_i}$, negative inverse, or $h(\eta_i) = \exp(\eta_i)$ and $g(\mu_i) = \ln(\mu_i)$, log-link.

Requirements

- one-to-one (inverse exists)
 - twice differential (for score function and expected Fisher information matrix)
-

Canonical link

$$\eta_i = \theta_i$$

so

$$g(\mu_i) = \theta_i$$

When the canonical link is used some of the results for the GLM (to be studied in the next sections) are simplified.

Likelihood inference set-up

$$\theta_i \leftrightarrow \mu_i \leftrightarrow \eta_i \leftrightarrow \beta$$

A more informative drawing made in class.

See class notes or Fahrmeir et al (2015), Section 5.8.2 for the derivation of the loglikelihood, score and expected Fisher information matrix.

Loglikelihood

$$l(\beta) = \sum_{i=1}^n l_i(\beta) = \sum_{i=1}^n \frac{1}{\phi} (y_i \theta_i - b(\theta_i)) w_i + \sum_{i=1}^n c(y_i, \phi, w_i)$$

Remark: the part of the loglikelihood involving both the data and the parameter of interest is for a *canonical link* equal to

$$\sum_{i=1}^n y_i \theta_i = \sum_{i=1}^n y_i \mathbf{x}_i^T \beta = \sum_{i=1}^n y_i \sum_{j=1}^p x_{ij} \beta_j = \sum_{j=1}^p \beta_j \sum_{i=1}^n y_i x_{ij}$$

Score function

$$s(\beta) = \sum_{i=1}^n \frac{(y_i - \mu_i) \mathbf{x}_i h'(\eta_i)}{\text{Var}(Y_i)} = \mathbf{X}^T \mathbf{D} \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu})$$

where $\Sigma = \text{diag}(\text{Var}(Y_i))$ and $\mathbf{D} = \text{diag}(h'(\eta_i))$ (derivative wrt η_i).

Remark: observe that $s(\beta) = 0$ only depends on the distribution of Y_i through μ_i and $\text{Var}(Y_i)$.

Canonical link:

$$s(\beta) = \sum_{i=1}^n \frac{(y_i - \mu_i) \mathbf{x}_i w_i}{\phi}$$

since $\frac{\partial \mu_i}{\partial \eta_i} = b'(\theta_i)$.

Expected Fisher information matrix for the GLM and covariance for $\hat{\beta}$

$$F_{[h,l]}(\beta) = \sum_{i=1}^n \frac{x_{ih} x_{il} (h'(\eta_i))^2}{\text{Var}(Y_i)}$$

$$F(\beta) = \mathbf{X}^T \mathbf{W} \mathbf{X}$$

where $\mathbf{W} = \text{diag}(\frac{h'(\eta_i)^2}{\text{Var}(Y_i)})$.

Canonical link:

$$\frac{\partial^2 l_i}{\partial \beta_j \partial \beta_l} = -\frac{x_{ij} w_i}{\phi} \left(\frac{\partial \mu_i}{\partial \beta_l} \right)$$

which do not contain any random variables, so the observed must be equal to the expected Fisher information matrix.

Fisher scoring and iterated reweighted least squares (IRWLS)

Details on the derivation: IRWLS

$$\beta^{(t+1)} = \beta^{(t)} + F(\beta^{(t)})^{-1} s(\beta^{(t)})$$

Insert formulas for expected Fisher information and score function.

$$\beta^{(t+1)} = (\mathbf{X}^T \mathbf{W}(\beta^{(t)}) \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}(\beta^{(t)}) \tilde{\mathbf{y}}_i^{(t)}$$

where \mathbf{W} is as before $\mathbf{W} = \text{diag}(\frac{h'(\eta_i)^2}{\text{Var}(Y_i)})$ - but now the current version of $\beta^{(t)}$ is used. The diagonal elements are called the *working weights*. The $\tilde{\mathbf{y}}_i^{(t)}$ is called the *working response vector* and has element i given as

$$\hat{y}_i^{(t)} = \mathbf{x}_i^T \beta^{(t)} + \frac{y_i - h(\mathbf{x}_i^T \beta^{(t)})}{h'(\mathbf{x}_i^T \beta^{(t)})}.$$

Remark: Convergens? With full rank of \mathbf{X} and positive diagonal elements of \mathbf{W} we are certain that the inverse will exist, but there might be that the temporary version of \mathbf{W} can cause problems.

See what is output from `glm`- observe working weights as `weights`..

```
fitgrouped = glm(cbind(y, n - y) ~ ldose, family = "binomial", data = investr::beetle)
names(fitgrouped)
fitgrouped$weights
fitgrouped$residuals
```

```
## [1] "coefficients"      "residuals"          "fitted.values"
## [4] "effects"            "R"                  "rank"
## [7] "qr"                 "family"             "linear.predictors"
## [10] "deviance"           "aic"                "null.deviance"
## [13] "iter"               "weights"            "prior.weights"
## [16] "df.residual"        "df.null"            "y"
## [19] "converged"          "boundary"           "model"
## [22] "call"               "formula"            "terms"
## [25] "data"               "offset"             "control"
## [28] "method"             "contrasts"          "xlevels"
##      1      2      3      4      5      6      7
## 3.254867 8.227383 14.321313 13.378893 10.261055 5.156671 2.653398
##      8
## 1.230713
##      1      2      3      4      5      6
## 0.78115418 0.38388091 -0.31082206 -0.44081641 0.18557365 -0.05641516
##      7      8
## 0.67002811 1.02139898
```

Estimator for dispersion parameter

Let data be grouped as much as possible. With G groups (covariate pattern) with n_i observations for each group (then $n = \sum^G n_i = n$):

$$\hat{\phi} = \frac{1}{G - p} \sum_{i=1}^G \frac{(y_i - \hat{\mu}_i)^2}{b''(\theta_i)/w_i}$$

The motivation behind this estimator is as follows:

$$\text{Var}(Y_i) = \phi b''(\theta_i)/w_i \Leftrightarrow \phi = \text{Var}(Y_i)/(b''(\theta_i)/w_i)$$

Distribution of the MLE

As before we have that maximum likelihood estimator $\hat{\beta}$ asymptotically follows the multivariate normal distribution with mean β and covariance matrix equal to the inverse of the expected Fisher information

matrix. This is also true when we replace the unknown β with the estimated $\hat{\beta}$ for the expected Fisher information matrix.

$$\hat{\beta} \approx N_p(\beta, F^{-1}(\hat{\beta}))$$

and with

$$F(\hat{\beta}) = \mathbf{X}^T \hat{\mathbf{W}} \mathbf{X}$$

where $\hat{\mathbf{W}}$ denotes that $\hat{\beta}$ is used then calculating $\mathbf{W} = \text{diag}(\frac{h'(\eta_i)^2}{\text{Var}(Y_i)})$.

What about the distribution of $\hat{\beta}, \hat{\phi}$?

The concept of orthogonal parameters

Hypothesis testing

Same as before - for the Wald we insert the formula for the covariance matrix of $\hat{\beta}$, for the LRT we insert the loglikelihoods and for the score test we insert formulas for the score function and expected Fisher information matrix.

Model assessment and model choice

Pearson and deviance statistic

Group observations together in groups of maximal size (covariate patterns? interval versions thereof?). Group i has n_i observations, and there are G groups. Asymptotic distribution correct if all groups have big n_i . For individual data asymptotic results can not be trusted.

Deviance

$$D = -2 \left[\sum_{i=1}^g (l_i(\hat{\mu}_i) - l_i(\bar{y}_i)) \right]$$

with approximate χ^2 -distribution with $G - p$ degrees of freedom.

Pearson:

$$X_P^2 = \sum_{i=1}^G \frac{(y_i - \hat{\mu}_i)^2}{v(\hat{\mu}_i)/w_i}$$

with approximate $\phi \cdot \chi^2$ -distribution with $G - p$ degrees of freedom.

Remember that the variance function $v(\hat{\mu}_i) = b''(\theta_i)$ (this is a function of μ_i because $\mu_i = b'(\theta_i)$).

AIC

Let p be the number of regression parameters in our model.

$$\text{AIC} = -2 \cdot l(\hat{\beta}) + 2p$$

If the dispersion parameter is estimated use $(p + 1)$ in place of p .

Interactive session

Work with Problem 1 and 2 in IL, and work on Problems 3-5 by yourself.

If you have more time after Problem 1 and 2, look through the theoretical proofs and derivations listed under “Learning material” on the top of this Module page.

Problem 1: Exam 2011, problem 3

a) Define the class of generalized linear models (GLMs), and explicitly list all requirements for each part of the model.

b) Below are three likelihoods, three link functions, and three linear components listed. Explain which *combinations* that give valid GLMs (8 in total), and also comment on these models (you do not have to mathematically prove which are valid).

Likelihoods:

1. Gaussian, $Y \sim N(\mu, \sigma^2)$
2. Binomial, $Y \sim Bin(n, \pi)$, where n is not fixed (hence is unknown and be estimated)
3. Poisson, $Y \sim Poisson(\lambda)$

Link functions:

1. $\eta = \cos(\mu)$
2. $\eta = \mu$
3. $\eta = \log(\mu)$

Linear components:

1. $\eta = \beta_0 + \beta_1 x_1 + \beta_2 x_2$
 2. $\eta = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2$
 3. $\eta = \beta_0 + \beta_1 x_1 + \beta_1^2 x_2$
-

Problem 2: December 2005, Problem 2 (modified)

1. Derive the formula for the (scaled) deviance for the binomial distribution.
 2. The covariance matrix for the estimated coefficients are given as $\text{Cov}(\hat{\beta}) = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1}$ where \mathbf{X} is the design matrix.
 - a. (New) The matrix \mathbf{W} is a diagonal matrix. What is on the diagonal?
 - b. Calculate the elements of \mathbf{W} for a Poisson regression- both with log and identity link. Compare.
 - c. Calculate the elements of \mathbf{W} for a binary regression - both with logit and identity link. Compare.
 - d. (New) Which insight did this give you into the role of the link function and its effect on the covariance for the parameter estimates?
-

Problem 3: Exam 2006, problem 2 (a, b, d)

Let Y_1, Y_2, \dots, Y_N be independent and exponentially distributed random variables, where Y_i has the density

$$f(y_i; \alpha_i) = \alpha_i e^{-\alpha_i y_i} \text{ for } y_i > 0, \alpha_i > 0, i = 1, 2, \dots, N.$$

- a) Show that the distribution of Y_i comes from the exponential family. Use the general formulas to find $E(Y_i)$ and $\text{Var}(Y_i)$ as functions of α_i .
- b) Show that the log-likelihood for the data y_1, \dots, y_n can be written as

$$l = \sum_{i=1}^N \{-\alpha_i y_i + \ln \alpha_i\}$$

Use this to show that the deviance for a generalized linear model with estimated expectations $\hat{\mu}_i = \hat{y}_i$ is

$$D = 2 \sum_{i=1}^N \left\{ \frac{y_i - \hat{y}_i}{\hat{y}_i} - \ln \left(\frac{y_i}{\hat{y}_i} \right) \right\}$$

- d) We want to test the null hypothesis

$$H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_N = \alpha$$

against the alternative hypothesis that at least one α_i is different from the others.

Use b) to find a test statistic for this problem.

What distribution does this test statistic have under H_0 ?

New: could you also use the Wald test? Explain how to do that (no need to calculate the mathematically explicit test statistic).

Problem 4: Exam 2007, problem 2 a, b, c

Assume that Y_1, \dots, Y_N are independent continuous distributed random variables, where the density of Y_i is given by

$$f(y_i; \gamma_i) = \begin{cases} \frac{\gamma_i^2}{2} y_i e^{-\gamma_i y_i} & \text{for } y_i \geq 0 \\ 0 & \text{else} \end{cases}$$

where γ_i is a scalar parameter.

- a) Show that the distribution of Y_i comes from the exponential family. Hint: usually we choose to let $\phi = \frac{1}{2}$. Use the general formulas to show that $E(Y_i) = 2/\gamma_i$ and $\text{Var}(Y_i) = 2/\gamma_i^2$.

Assume a GLM for Y_1, \dots, Y_N where the distribution of Y_i is as above for all i , with the following link function:

$$\eta = g(\mu) = \ln(\mu) = x^T \beta$$

where $x, \beta \in \mathbb{R}^p$ and μ is the expected value of y .

- b) Use general formulas to find the score vector $s(\beta) = [s_1(\beta), \dots, s_p(\beta)]^T$ and the expected Fisher information matrix $F(\beta) = [F_{ij}(\beta)]_{i,j=1}^p$, expressed using $y_1, \dots, y_N, \beta, N$ and the covariates x . Write down the equation that can be used to find the MLE for β . Note that this is a recursive equation.

c) Write down the log-likelihood for the model above. Use this to find the deviance D for the model as a function of y_1, \dots, y_N and $\hat{y}_1, \dots, \hat{y}_N$, where \hat{y}_i is the estimated expected value of y_i . Find an expression for the deviance residuals d_i using y_i and \hat{y}_i .

Problem 5: Exam UiO December 2017, Problem 2

We assume that the random variable Λ is gamma distributed with pdf

$$f(\lambda; \nu, \mu) = \frac{(\nu/\mu)^\nu}{\Gamma(\nu)} \lambda^{\nu-1} e^{-\lambda/\mu}; \lambda > 0$$

and further that given $\Lambda = \lambda$, the random variable Y is Poisson distributed with parameter λ . Thus the conditional pmf of Y given $\Lambda = \lambda$ takes the form

$$P(Y = y|\lambda) = \frac{\lambda^y}{y!} \exp(-\lambda), \quad y = 0, 1, 2, \dots$$

a) Show that the marginal pmf of Y is given by

$$p(y; \mu, \nu) = \frac{\Gamma(y + \nu)}{\Gamma(\nu)\Gamma(y + 1)} \left(\frac{\mu}{\mu + \nu}\right)^y \left(\frac{\nu}{\mu + \nu}\right)^\nu; \quad y = 0, 1, 2, \dots$$

This is the negative binomial distribution.

We then assume that the parameter ν is fixed, and consider the random variable $Y^* = Y/\nu$. Note that

$$P(Y^* = y^*) = P(Y = ky^*) \text{ for } y^* = 0, \frac{1}{k}, \frac{2}{k}, \dots$$

so Y^* has pmf

$$p^*(y^*; \mu, \nu) = \frac{\Gamma(\nu y^* + \nu)}{\Gamma(\nu)\Gamma(\nu y^* + 1)} \left(\frac{\mu}{\mu + \nu}\right)^{\nu y^*} \left(\frac{\nu}{\mu + \nu}\right)^\nu; \quad y^* = 0, \frac{1}{k}, \frac{2}{k}, \dots$$

b) Show that the pmf of Y^* is an exponential family

$$\exp \left\{ \frac{y\theta - b(\theta)}{\phi} w + c(y, \phi, w) \right\},$$

with $\theta = \log(\mu/(\mu + \nu))$, $b(\theta) = -\log(1 - e^\theta)$, $w = 1$ and $\phi = 1/\nu$

c) Use the expressions for $b(\theta)$ and ϕ to determine $E(Y^*)$ and $\text{Var}(Y^*)$. Show that $E(Y) = \mu$ and find $\text{Var}(Y)$.

Exam questions

December 2015

One of the important concepts we have discussed in this course, is deviance (for Gaussian regression, Poisson regression and logistic regression).

1. Explain what deviance is, and how it relates to residual sum of squares (RSS) for Gaussian regression. Remark 2017/2018: we have called this “sums of squares of errors - SSE”

2. Discuss how it relates to a likelihood ratio test (LRT) for comparing two nested regression models.
3. Discuss how deviance can be used to construct “ANOVA” tables for Poisson regression and logistic regression. Remark 2017/2018: these are called analysis of deviance tables.
4. Discuss how deviance can be used to define residuals, for Poisson regression and logistic regression.

Further reading

- A. Agresti (2015): “Foundations of Linear and Generalized Linear Models.” Wiley.