

Part II: Generalized Linear Models

Chapter II.2

Theory of Generalized Linear Models

Maximum Likelihood Estimation

Maximum Likelihood Estimation

Since GLMs assume a specific probability distribution for the responses that belongs to the exponential dispersion family (EDF), maximum likelihood estimation (MLE) procedures are used for parameter estimation, and a general formulation can be developed within the GLM set-up.

Log-Likelihood for GLMs' Parameters Vector β

Consider n independent responses $Y_i \sim \text{EDF}(\vartheta_i, \phi)$, with the associated function a possibly dependent on the observations, i.e. $a(\phi) = a(\phi; i)$, $i = 1, \dots, n$. Then the corresponding log-likelihood function (s. Definition II.2.3) is:

$$\ell(\beta) = \sum_{i=1}^n \ell_i = \sum_{i=1}^n \log f(y_i; \vartheta_i, \phi) = \sum_{i=1}^n \frac{y_i \vartheta_i - b(\vartheta_i)}{a(\phi; i)} + \sum_{i=1}^n c(y_i, \phi) .$$

It is a function of β due to $g(\mu) = \eta = \mathbf{X}\beta$ and $\mu = E(Y) = b'(\vartheta)$.

The first derivative of ℓ is the *score function* (s. Definition II.1.8):

$$S(\beta) = \frac{\partial \ell(\beta)}{\partial \beta} = \left(\frac{\partial \ell(\beta)}{\partial \beta_1}, \dots, \frac{\partial \ell(\beta)}{\partial \beta_p} \right)',$$

where $\frac{\partial \ell(\beta)}{\partial \beta_j} = \sum_{i=1}^n \frac{\partial \ell_i}{\partial \beta_j}$, $j = 1, \dots, p$. By the chain rule we further have

$$\frac{\partial \ell_i}{\partial \beta_j} = \frac{\partial \ell_i}{\partial \vartheta_i} \cdot \frac{\partial \vartheta_i}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_j} .$$

Derivation of $\frac{\partial \ell_i}{\partial \beta_j}$, $i=1, \dots, n$, $j=1, \dots, p$

$$\eta_i = g(\mu_i) = \sum_{k=1}^p \beta_k x_{ik}$$

$$\frac{\partial \ell_i}{\partial \beta_j} = \frac{\partial \ell_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_j} \quad \boxed{\text{**}}$$

$$(1) \frac{\partial \ell_i}{\partial \eta_i} = \frac{\partial}{\partial \eta_i} \left(\frac{y_i \eta_i - b(\eta_i)}{\alpha(\phi; i)} \right) = \frac{y_i - b'(\eta_i)}{\alpha(\phi; i)} = \frac{y_i - \mu_i}{\alpha(\phi; i)}$$

$$(2) \frac{\partial \mu_i}{\partial \theta_i} = \frac{\partial (b'(\eta_i))}{\partial \theta_i} = b''(\eta_i)$$

$$(3) \frac{\partial \eta_i}{\partial \beta_j} \stackrel{*}{=} \frac{\partial}{\partial \beta_j} \left(\sum_{k=1}^p \beta_k x_{ik} \right) = x_{ij}$$

$$\boxed{\text{**}} \stackrel{(1)}{\Rightarrow} \frac{\partial \ell_i}{\partial \beta_j} =$$

$$\frac{y_i - \mu_i}{\alpha(\phi; i)} \cdot \frac{1}{b''(\eta_i)} \frac{\partial \mu_i}{\partial \eta_i} x_{ij} = \\ = \text{Var}(Y_i)$$

$$\frac{y_i - \mu_i}{\text{Var}(Y_i)} \cdot \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} x_{ij}$$



► II.2.16 Likelihood equations for parameters vector β

Equating the components of the score function $S(\beta)$ to zero, the *likelihood equations* are obtained

$$\sum_{i=1}^n \left(\frac{y_i - \mathbb{E}(Y_i)}{\text{Var}(Y_i)} \cdot \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \cdot x_{ik} \right) = 0 , \quad k = 1, \dots, p ,$$

where $\eta_i = \sum_{j=1}^p x_{ij} \beta_j = g(\mu_i)$, when the link function is g .

► II.2.17 Remark

The likelihood equations above have the matrix form

$$\mathbf{X}' \mathbf{D} \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0} ,$$

where $\mathbf{V} = \text{diag}(\text{Var}(Y_1), \dots, \text{Var}(Y_n))$, $\mathbf{D} = \text{diag}\left(\frac{\partial \mu_1}{\partial \eta_1}, \dots, \frac{\partial \mu_n}{\partial \eta_n}\right)$ and $\mathbf{0}$ is an n -dimensional column vector of 0s. Although β does not appear in the equation above, it is there implicitly through μ , since $\mu_i = g^{-1}\left(\sum_{j=1}^p \beta_j x_{ij}\right)$.

► II.2.18 Remark (the mean-variance relation in EDF)

The likelihood equations of a GLM depend on the distribution of Y_i only through $E(Y_i)$ and $\text{Var}(Y_i)$, $i = 1, \dots, n$.

Furthermore,

- $\text{Var}(Y_i)$ depends on $\mu_i = E(Y_i)$ through a functional form $\text{Var}(Y_i) = v(\mu_i)$.
For example,
 - $Y_i \sim \mathcal{P}(\mu_i)$: $\text{Var}(Y_i) = \mu_i$ (identity function),
 - $Y_i \sim \mathcal{N}(\mu_i, \sigma^2)$: $\text{Var}(Y_i) = \sigma^2$ (constant function),
 - $Y_i \sim \mathcal{B}(m_i, \pi_i)/m_i$ with $\mu_i = E(Y_i) = \pi_i$: $\text{Var}(Y_i) = \frac{\mu_i(1-\mu_i)}{m_i}$.
- In EDF, this mean-variance relation *characterizes* the distribution.
Example: if $\text{Var}(Y_i) = E(Y_i)$, then Y_i has to be Poisson distributed.

II.2.19 Likelihood equations for GLMs with canonical link

Consider n independent responses $Y_i \sim \text{EDF}(\vartheta_i, \phi)$ with $a(\phi) = a(\phi; i)$ and $\text{E}(Y_i) = \mu_i$ $i = 1, \dots, n$. For a GLM with *canonical link*, the likelihood equations II.2.16 are simplified to

$$\sum_{i=1}^n \frac{1}{a(\phi; i)} (y_i - \mu_i) x_{ik} = 0, \quad k = 1, \dots, p.$$

Proof

Due to the canonical link, we have $\eta_i = g(\mu_i) = \vartheta_i$ and hence $\mu_i = g^{-1}(\vartheta_i) = g^{-1}(\eta_i)$, $i = 1, \dots, n$. Thus,

$$\frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} = \frac{\partial g^{-1}(\vartheta_i)}{\partial \vartheta_i} = \frac{\partial \mu_i}{\partial \vartheta_i} = \frac{\partial b'(\vartheta_i)}{\partial \vartheta_i} = b''(\vartheta_i), \quad i = 1, \dots, n.$$

Furthermore, it holds $\text{Var}(Y_i) = a(\phi; i)b''(\vartheta_i)$ (see Proposition II.2.4).

The result is derived by substituting $\text{Var}(Y_i)$ and $\frac{\partial g^{-1}(\eta_i)}{\partial \eta_i}$ in the likelihood equations II.2.16 by the expressions above.

Fisher Information Matrix: The expected Fisher information matrix has entries

$$E\left(-\frac{\partial^2 l(\beta)}{\partial \beta_j \partial \beta_k}\right) = E\left(\frac{\partial l(\beta)}{\partial \beta_j} \cdot \frac{\partial l(\beta)}{\partial \beta_k}\right), \quad j, k = 1, \dots, p$$

it holds for EDF (since the regularity conditions hold)

Since

$$l(\beta) = \sum_{i=1}^n l_i(\beta), \text{ we work on the } i\text{-th term of the}$$

log-likelihood and have

$$\begin{aligned} E\left(\frac{\partial l_i(\beta)}{\partial \beta_j} \frac{\partial l_i(\beta)}{\partial \beta_k}\right) &= E\left(\underbrace{\frac{(Y_i - \mu_i)}{\text{Var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_{ij}}}_{x_{ij}} \underbrace{x_{ik}}_{\frac{(Y_i - \mu_i)}{\text{Var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_{ik}}}\right) \\ &= x_{ij} x_{ik} \frac{(\partial \mu_i)^2}{[\text{Var}(Y_i)]^2} \underbrace{\frac{E(Y_i - \mu_i)^2}{\text{Var}(Y_i)}}_{= \text{Var } Y_i} \end{aligned}$$



Hence :

$$E\left(\frac{\partial \ell_i(\beta)}{\partial \beta_j} \frac{\partial \ell_i(\beta)}{\partial \beta_k}\right) = X_{ij} X_{ik} \left(\frac{\partial \mu_i}{\partial \eta_i}\right)^2 \cdot \frac{1}{\text{Var}(Y_i)}$$

and

$$E\left(\frac{\partial \ell(\beta)}{\partial \beta_j} \frac{\partial \ell(\beta)}{\partial \beta_k}\right) = \sum_{i=1}^n \frac{X_{ij} X_{ik}}{\text{Var}(Y_i)} \left(\frac{\partial \mu_i}{\partial \eta_i}\right)^2$$

Overall, for $W_{n \times n} = \text{diag} \left(\frac{(\partial \mu_i / \partial \eta_i)^2}{\text{Var}(Y_i)} \right)$

we get

$$I_F = X' W X$$

The s.e. of $\hat{\beta}_j$, $j=1, \dots, p$,
are the square roots
of the diagonal elements
of I_F^{-1} .

► II.2.20 Remark (likelihood equations for Poisson GLMs)

The Poisson GLM is the ‘standard’ GLM for count data (s. Chapters II.5 and II.6). It assumes that the random component is Poisson distributed and adopts the associated canonical link, i.e. the log-link $\log(\mu_i)$ (s. Example II.2.15). In this case, considering n independent responses $Y_i \sim \mathcal{P}(\mu_i)$, $i = 1, \dots, n$, the likelihood equations in II.2.19 become

$$\sum_{i=1}^n (y_i - \mu_i)x_{ik} = 0 , \quad k = 1, \dots, p ,$$

or, in matrix notation,

$$\mathbf{X}'(\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0} , \quad \mathbf{0} : n\text{-dimensional column vector of 0s} ,$$

since $a(\phi; i) = \phi/w_i$ with $\phi = 1$, $w_i = 1$ (s. Remark II.2.8) and $E(Y_i) = \mu_i$, $i = 1, \dots, n$.

Note that $\mathbf{X}'\mathbf{Y}$ is a *sufficient* statistic for $\boldsymbol{\beta}$ and the likelihood equations equate every β_k 's sufficient statistic value $\sum_{i=1}^n x_{ik}y_i$ to its expected value.

 If all explanatory variables are categorical, then the model is the *Poisson log-linear model* (s. Chapter II.5) while in the presence of continuous explanatory variables, the model is the so called *Poisson regression* (s. Chapter II.6).

► II.2.21 Remark (likelihood equations for binomial logit models)

If the responses to be modeled are proportions of success, then assuming a sample of n independent random proportions, i.e. $Y_i \sim \mathcal{B}(m_i, \pi_i)/m_i$, $i = 1, \dots, n$, the associate canonical link is the *logit* link $\log(\frac{\pi_i}{1-\pi_i})$ (s. Example II.2.15) and the likelihood equations in II.2.19 become

$$\sum_{i=1}^n m_i(y_i - \pi_i)x_{ik} = 0 , \quad k = 1, \dots, p ,$$

since $a(\phi; i) = \phi/w_i$ with $\phi = 1$, $w_i = m_i$ (s. Remark II.2.8) and $E(Y_i) = \mu_i = \pi_i$, $i = 1, \dots, n$.

► II.2.22 Theorem (large sample Normal distribution of MLE $\hat{\beta}$)

Under standard *regularity conditions*^a, the maximum likelihood estimator (MLE) of β , $\hat{\beta}$, has an approximate normal distribution. Thus, for large n it holds

$$\hat{\beta} \sim \mathcal{N}_p \left(\beta, \text{Cov}(\hat{\beta}) \right).$$

The **asymptotic covariance matrix** of $\hat{\beta}$ is $\text{Cov}(\hat{\beta}) = \mathcal{I}_F^{-1}$, where \mathcal{I}_F is the *expected Fisher information matrix*^b

$$\mathcal{I}_F = \mathbb{E} \left(\frac{\partial \ell}{\partial \beta} \frac{\partial \ell}{\partial \beta'} \right) = \mathbb{E} \left(-\frac{\partial^2 \ell}{\partial \beta \partial \beta'} \right) = \mathbf{X}' \mathbf{W} \mathbf{X},$$

with $\mathbf{W} = \text{diag} \left(\frac{(\partial \mu_1 / \partial \eta_1)^2}{\text{Var}(Y_1)}, \dots, \frac{(\partial \mu_n / \partial \eta_n)^2}{\text{Var}(Y_n)} \right)$.

^aThey mainly require that β is in the interior of the parameter space and has fixed dimension as n increases (s. II.1.15 or Cox & Hinkley (2000, *Theoretical Statistics*, Chapman & Hall/CRC, Section 9.1).

^bsee Definition II.1.13

The asymptotic covariance matrix is estimated as

$$\widehat{\text{Var}}(\hat{\beta}) = \mathcal{I}_F^{-1}(\hat{\beta}) = \left(\mathbf{X}' \widehat{\mathbf{W}} \mathbf{X} \right)^{-1},$$

where $\widehat{\mathbf{W}}$ is \mathbf{W} evaluated at $\hat{\beta}$.

II.2.22 Theorem (large sample Normal distribution of MLE $\hat{\beta}$)

Under standard *regularity conditions*^a, the maximum likelihood estimator (MLE) of β , $\hat{\beta}$, has an approximate normal distribution. Thus, for large n it holds

$$\hat{\beta} \sim \mathcal{N}_p(\beta, \text{Cov}(\hat{\beta})).$$

The **asymptotic covariance matrix** of $\hat{\beta}$ is $\text{Cov}(\hat{\beta}) = \mathcal{I}_F^{-1}$, where \mathcal{I}_F is the *expected Fisher information matrix*^b

$$\mathcal{I}_F = E\left(\frac{\partial \ell}{\partial \beta} \frac{\partial \ell}{\partial \beta'}\right) = E\left(-\frac{\partial^2 \ell}{\partial \beta \partial \beta'}\right) = \mathbf{X}' \mathbf{W} \mathbf{X},$$

with $\mathbf{W} = \text{diag}\left(\frac{(\partial \mu_1 / \partial \eta_1)^2}{\text{Var}(Y_1)}, \dots, \frac{(\partial \mu_n / \partial \eta_n)^2}{\text{Var}(Y_n)}\right)$.

\mathbf{W} (and \mathcal{I}_F) depend on the link function!

^aThey mainly require that β is in the interior of the parameter space and has fixed dimension as n increases (s. II.1.15 or Cox & Hinkley (2000, *Theoretical Statistics*, Chapman & Hall/CRC, Section 9.1).

^bsee Definition II.1.13

The asymptotic covariance matrix is estimated as

$$\widehat{\text{Cov}}(\hat{\beta}) = \mathcal{I}_F^{-1}(\hat{\beta}) = (\mathbf{X}' \widehat{\mathbf{W}} \mathbf{X})^{-1},$$

where $\widehat{\mathbf{W}}$ is \mathbf{W} evaluated at $\hat{\beta}$.

► II.2.23 Definition (observed information matrix)

Let $\ell(\beta)$ be the log-likelihood of a GLM with associated parameter vector $\beta = (\beta_1, \dots, \beta_p)'$. The observed information matrix is the following $p \times p$ matrix

$$\mathcal{I}_F^{obs} = \left(-\frac{\partial^2 \ell}{\partial \beta_k \partial \beta_r} \right) = -\mathcal{H},$$

where \mathcal{H} is known as the *Hessian matrix*. It holds:

$$\mathbf{W}_c = \text{diag}(2 \ 8 \ 18)$$

$$-\mathbf{H} = 20 \ -18$$

$$\mathcal{I}_F = \mathbb{E} \left(\mathcal{I}_F^{obs} \right) = \mathbb{E} (-\mathcal{H}). \quad -18 \ 26$$

► II.2.24 Information matrix for GLMs with canonical link

For a GLM with canonical link function, since $\eta_i = \vartheta_i$, it follows that (s. Proof of II.2.19)

$$\frac{\partial \mu_i}{\partial \eta_i} = \frac{\partial b'(\vartheta_i)}{\partial \vartheta_i} = b''(\vartheta_i), \quad i = 1, \dots, n, \text{ leading to}$$

$$\mathcal{H} = -\mathbf{X}' \mathbf{W}_c \mathbf{X},$$

where $\mathbf{W}_c = \text{diag}(w_1, \dots, w_n)$ with $w_i = \frac{b''(\vartheta_i)}{a(\phi; i)}$, **independent** of \mathbf{y} . Hence

$$\boxed{\mathcal{I}_F = \mathbb{E} (-\mathcal{H}) = -\mathcal{H} = \mathcal{I}_F^{obs}}$$

Existence of MLEs

II.2.25 Remark

A necessary **but** not sufficient condition for the existence and uniqueness of the MLEs is that the model matrix \mathbf{X} is of full rank.

For many GLMs with full-rank \mathbf{X} (including Poisson log-linear and binomial logit models), the Hessian matrix is negative definite and the log-likelihood strictly concave. In such cases, the MLEs of the model parameters exist uniquely under quite general conditions ^a.

^aWedderburn (1976). On the existence and uniqueness of the maximum likelihood estimates for certain generalized linear models, *Biometrika*, 63, 27-32.

Fitting GLMs

The likelihood equations of a GLM are usually nonlinear in β and do not lead to closed form expressions for the MLE $\hat{\beta}$. The two methods usually applied for solving the likelihood equations iteratively are

- ① the *Newton-Raphson*, and
- ② the *Fisher scoring*.

► II.2.26 Algorithm (Newton-Raphson)

The Newton-Raphson is an iterative algorithm having following steps, adjusted in our estimation problem of β :

- ① Initial guess ($t = 0$): Set $\beta^{(t)} = \beta^{(0)}$.
- ② For stage $t + 1$, $t = 0, 1, 2, \dots$:
 - If $\beta^{(t)}$ is the value assigned to $\hat{\beta}$ at the previous stage t of the iterative procedure, obtain the next guess $\beta^{(t+1)}$ through the updating equations of the Newton-Raphson algorithm

$$\beta^{(t+1)} = \beta^{(t)} - (\mathcal{H}^{(t)})^{-1} S(\beta^{(t)}), \quad (\text{II.4})$$

where $S(\beta^{(t)})$ and $\mathcal{H}^{(t)}$ are the score function $S(\beta)$ and the Hessian matrix \mathcal{H} evaluated at $\beta^{(t)}$. For matrix inversion to be possible, $\mathcal{H}^{(t)}$ has to be non-singular.

- Check for convergence: The algorithm converges and stops, say after t_c iterations, when a termination criterion is met, leading to $\hat{\beta} = \beta^{(t_c)}$. A termination criterion checks whether $\beta^{(t)}$ and $\beta^{(t+1)}$ are sufficient close, e.g. $\sum_{j=1}^p |\beta_j^{t+1} - \beta_j^t| \leq c$, for some prespecified small c .

► II.2.27 Remark

The updating equation (II.4) is derived by determining the point $\hat{\beta}$ at which $\ell(\beta)$ is maximized, when $\ell(\beta)$ is approximated near $\beta^{(t)}$ by the terms up to the second order in the Taylor expansion.

Newton-Raphson Method : quadratic (2nd order) approximation

Consider a 2nd order Taylor expansion of f around a (univariate):

$$f(x) \approx f(a) + f'(a)(x-a) + \frac{1}{2}f''(a)(x-a)^2, \quad x, a \in \mathbb{R}$$

For $x, a \in \mathbb{R}^P$:

$$f(x) \approx f(a) + S(a)'(x-a) + \frac{1}{2}(x-a)'H(a)(x-a),$$

where $S(a) = \left(\frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_p} \right) \Big|_{x=a}$ and

$$H(a) = \left(\frac{\partial^2 f(x)}{\partial x_j \partial x_k} \right) \Big|_{x=a}.$$

In our case f is the log-likelih.
 β , x is the parameter β
and $a = \beta^{(t)}$. Thus $S(\beta^{(t)})$
and $H(\beta^{(t)})$ are the score funct.
and Hessian matrix evaluated at $\beta^{(t)}$.



Thus:

$$l(\beta) = Q(\beta^{(t)}) + S(\beta^{(t)})'(\underline{\beta} - \underline{\beta}^{(t)}) + \frac{1}{2} (\underline{\beta} - \underline{\beta}^{(t)})' J L(\beta^{(t)}) (\underline{\beta} - \underline{\beta}^{(t)})$$

$$\Rightarrow \frac{\partial l(\beta)}{\partial \beta} = S(\beta^{(t)}) + J L(\beta^{(t)}) (\underline{\beta} - \underline{\beta}^{(t)})$$

Solving $\frac{\partial l(\beta)}{\partial \beta} = 0$ for non-singular JL we
get

$$\underline{\beta} = \underline{\beta}^{(t)} - (J L(\beta^{(t)}))^{-1} S(\beta^{(t)})$$

The convergence of $\beta^{(t)}$ to $\hat{\beta}$ for the N-R method is usually fast. For large t :

$$|\beta_j^{(t+1)} - \hat{\beta}_j| \leq c^* |\beta_j^{(t)} - \hat{\beta}_j|^2, \text{ some } c^*, \forall j$$

"2nd order convergence"



► II.2.28 Algorithm (Fisher Scoring)

The *Fisher scoring* is similar to the *Newton-Raphson* algorithm with the only difference being that it is based on the **expected information matrix** \mathcal{I}_F , instead of the observed information matrix $\mathcal{I}_F^{obs} = -\mathcal{H}$.

In particular, the updating equations for the Fisher scoring algorithm are

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} + \left(\mathcal{I}_F^{(t)} \right)^{-1} S(\boldsymbol{\beta}^{(t)}) ,$$

where $\mathcal{I}_F^{(t)}$ is \mathcal{I}_F evaluated at $\boldsymbol{\beta}^{(t)}$.

► II.2.29 Remark

For GLMs with *canonical link*, the Newton-Raphson and Fischer scoring algorithms coincide, since the observed and expected information matrices are equal (see Remark II.2.24).



II.2.28 Algorithm (Fisher Scoring)

The *Fisher scoring* is similar to the *Newton-Raphson* algorithm with the only difference being that it is based on the **expected information matrix** \mathcal{I}_F , instead of the observed information matrix $\mathcal{I}_F^{obs} = -\mathcal{H}$.

In particular, the updating equations for the Fisher scoring algorithm are

$$\beta^{t+1} = \beta^t + \left(\mathcal{I}_F^{(t)} \right)^{-1} S(\beta^{(t)}) ,$$

where $\mathcal{I}_F^{(t)}$ is \mathcal{I}_F evaluated at $\beta^{(t)}$.

II.2.29 Remark

For GLMs with *canonical link*, the Newton-Raphson and Fischer scoring algorithms coincide, since the observed and expected information matrices are equal (see Remark II.2.24).

Advantage of Fisher Scoring over N-R: The estimated asymptotic cov. matrix $\hat{\mathcal{I}}_F^{-1}$ of $\hat{\beta}$ occurs as a byproduct of this algorithm: $\text{Cov}(\hat{\beta}) = (\hat{X}' \hat{W} \hat{X})^{-1} = \hat{\mathcal{I}}_F^{-1}$

► II.2.30 Remark (initial values)

The data \mathbf{y} can be used as the initial estimates of $\boldsymbol{\mu}$ (for both algorithms). This determines the first estimate of \mathbf{W} and hence $\boldsymbol{\beta}$.

► II.2.31 Corollary (covariance matrix of fitted values)

For large sample sizes, since $\hat{\boldsymbol{\eta}} = g(\hat{\boldsymbol{\mu}}) = \mathbf{X}\hat{\boldsymbol{\beta}}$ and $\text{Cov}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}$ (see Theorem II.2.22), it holds

$$\text{Cov}(\hat{\boldsymbol{\eta}}) = \mathbf{X}\text{Cov}(\hat{\boldsymbol{\beta}})\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'.$$

Furthermore, since $\hat{\boldsymbol{\mu}} = g^{-1}(\hat{\boldsymbol{\eta}})$, the asymptotic covariance matrix $\text{Cov}(\hat{\boldsymbol{\mu}})$ can be obtained from $\text{Cov}(\hat{\boldsymbol{\eta}})$ by the *delta method*

$$\text{Cov}(\hat{\boldsymbol{\mu}}) = \left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} \right) \text{Cov}(\hat{\boldsymbol{\eta}}) \left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} \right)' = \mathbf{D} \left(\mathbf{X}(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}' \right) \mathbf{D},$$

where $\mathbf{D} = \text{diag} \left(\frac{\partial \mu_1}{\partial \eta_1}, \dots, \frac{\partial \mu_n}{\partial \eta_n} \right)$ (see Remark II.2.17).



II.2.30 Remark (initial values)

The data \mathbf{y} can be used as the initial estimates of $\boldsymbol{\mu}$ (for both algorithms). This determines the first estimate of \mathbf{W} and hence $\boldsymbol{\beta}$.

II.2.31 Corollary (covariance matrix of fitted values)

For large sample sizes, since $\hat{\boldsymbol{\eta}} = g(\hat{\boldsymbol{\mu}}) = \mathbf{X}\hat{\boldsymbol{\beta}}$ and $\text{Cov}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}$ (see Theorem II.2.22), it holds

$$\text{Cov}(\hat{\boldsymbol{\eta}}) = \mathbf{X}\text{Cov}(\hat{\boldsymbol{\beta}})\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'.$$

Furthermore, since $\hat{\boldsymbol{\mu}} = g^{-1}(\hat{\boldsymbol{\eta}})$, the asymptotic covariance matrix $\text{Cov}(\hat{\boldsymbol{\mu}})$ can be obtained from $\text{Cov}(\hat{\boldsymbol{\eta}})$ by the *delta method*

$$\text{Cov}(\hat{\boldsymbol{\mu}}) = \left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} \right) \text{Cov}(\hat{\boldsymbol{\eta}}) \left(\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} \right)' = \mathbf{D} \left(\mathbf{X}(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}' \right) \mathbf{D}',$$

where $\mathbf{D} = \text{diag} \left(\frac{\partial \mu_1}{\partial \eta_1}, \dots, \frac{\partial \mu_n}{\partial \eta_n} \right)$ (see Remark II.2.17).