

Advanced Statistical Modelling III (Epiphany term)

Department of Mathematical Sciences at Durham University

09 Jan, 2025

Contents

General Information	5
1 Introduction and Review	7
1.1 Introduction	7
1.2 Random Vectors and Random Matrices: A Review	7
1.3 Generalised Linear Models: A Review	9
1.4 Exercises	10
2 Estimation	13
2.1 Likelihood Function	13
2.2 Log-Likelihood Function	14
2.3 Score Function and Score Equation	14
2.4 Fisher Information	15
2.5 Example: Poisson Regression	16
2.6 Properties of $\mathbf{S}(\boldsymbol{\beta})$ and $\mathbf{F}(\boldsymbol{\beta})$	17
2.7 Matrix Notation	20
2.8 Iterative Solution of $\mathbf{S}(\hat{\boldsymbol{\beta}}) = 0$	21
2.9 Practical Example: US Polio Data	25
2.10 Estimation of Dispersion Parameter ϕ	31
2.11 Asymptotic Properties of $\hat{\boldsymbol{\beta}}$	33
2.12 Exercises	35

General Information

- These are the lecture notes for the second term of the module MATH3411 – Advanced Statistical Modelling III of Durham University’s degree for Mathematics and Statistics.
- **Acknowledgements:** This material is based on the lecture notes in previous modules taught by Dr Samuel Jackson, Dr Jochen Einbeck, Dr Ric Crossman, Dr Emmanuel Ogundimu, Dr Cuong Nguyen, Dr Ian Jermyn, Dr Louis Aslett, and Dr Reza Drikvandi.

Chapter 1

Introduction and Review

1.1 Introduction

In the first term, we have learned about categorical data analysis and some basics of generalised linear models (GLMs). In this term, we will continue to explore GLMs in more detail and study some of its more general variants. In particular, we will learn:

- How to estimate the parameters of a GLM from data.
- How to make a prediction and do inference once a GLM has been fitted.
- How to perform deviance analysis with GLMs.
- How to reduce overdispersion using quasi-likelihood methods.
- How to model repeated measures data using marginal models.
- How to model mixed effects using linear mixed models (LMMs) and generalised linear mixed models (GLMMs).

For the rest of this chapter, we will review some random vector and random matrix identities that will be useful later. Then we will review the basics of GLMs.

1.2 Random Vectors and Random Matrices: A Review

- A *random vector* is a vector of random variables:

$$\mathbf{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}.$$

- The mean or expectation of \mathbf{X} is defined as:

$$E[\mathbf{X}] = \begin{pmatrix} E[X_1] \\ \vdots \\ E[X_n] \end{pmatrix}.$$

- A *random matrix* is a matrix of random variables:

$$\mathbf{Z} = (Z_{ij}) = \begin{pmatrix} Z_{11} & \dots & Z_{1m} \\ \vdots & \ddots & \vdots \\ Z_{n1} & \dots & Z_{nm} \end{pmatrix}.$$

- The expectation of \mathbf{Z} is defined as:

$$E[\mathbf{Z}] = (E[Z_{ij}]) = \begin{pmatrix} E[Z_{11}] & \dots & E[Z_{1m}] \\ \vdots & \ddots & \vdots \\ E[Z_{n1}] & \dots & E[Z_{nm}] \end{pmatrix}.$$

- Some properties of random vectors and random matrices:
 - If \mathbf{a} is a constant (i.e., non-random) vector, $E[\mathbf{a}] = \mathbf{a}$.
 - If \mathbf{A} is a constant matrix, $E[\mathbf{A}] = \mathbf{A}$.
 - $E[\mathbf{X} + \mathbf{Y}] = E[\mathbf{X}] + E[\mathbf{Y}]$ for any random matrices \mathbf{X} and \mathbf{Y} .
 - $E[\mathbf{A}\mathbf{X}] = \mathbf{A}E[\mathbf{X}]$ for a constant matrix \mathbf{A} and a random matrix \mathbf{X} .
 - More generally, $E[\mathbf{A}\mathbf{X}\mathbf{B} + \mathbf{C}] = \mathbf{A}E[\mathbf{X}]\mathbf{B} + \mathbf{C}$ for constant matrices \mathbf{A} , \mathbf{B} and \mathbf{C} .
- Let \mathbf{X} be a random vector. The *covariance matrix* of \mathbf{X} is defined as:

$$\text{cov}(\mathbf{X}) = (\text{cov}(X_i, X_j)) = \begin{pmatrix} \text{var}(X_1) & \text{cov}(X_1, X_2) & \dots & \text{cov}(X_1, X_n) \\ \text{cov}(X_2, X_1) & \text{var}(X_2) & \dots & \text{cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_n, X_1) & \text{cov}(X_n, X_2) & \dots & \text{var}(X_n) \end{pmatrix}.$$

- We can also write:

$$\begin{aligned} \text{cov}(\mathbf{X}) &= E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{X} - E[\mathbf{X}])^T] \\ &= E \left[\begin{pmatrix} X_1 - E[X_1] \\ \vdots \\ X_n - E[X_n] \end{pmatrix} (X_1 - E[X_1], \dots, X_n - E[X_n]) \right]. \end{aligned}$$

- Some properties of covariance matrices:
 - They are symmetric: $\text{cov}(\mathbf{X}) = \text{cov}(\mathbf{X})^T$.
 - If \mathbf{a} is a constant vector, $\text{cov}(\mathbf{X} + \mathbf{a}) = \text{cov}(\mathbf{X})$.
 - If \mathbf{A} is a constant matrix, $\text{cov}(\mathbf{A}\mathbf{X}) = \mathbf{A}\text{cov}(\mathbf{X})\mathbf{A}^T$.
 - $\text{cov}(\mathbf{X}) = E[\mathbf{X}\mathbf{X}^T] - E[\mathbf{X}]E[\mathbf{X}]^T$.
- Let $\mathbf{X}_{n \times 1}$ and $\mathbf{Y}_{m \times 1}$ be random vectors where n and m could be different. The covariance matrix of \mathbf{X} and \mathbf{Y} is defined as:

$$\text{cov}(\mathbf{X}, \mathbf{Y}) = (\text{cov}(X_i, Y_j)) = \begin{pmatrix} \text{cov}(X_1, Y_1) & \text{cov}(X_1, Y_2) & \dots & \text{cov}(X_1, Y_m) \\ \text{cov}(X_2, Y_1) & \text{cov}(X_2, Y_2) & \dots & \text{cov}(X_2, Y_m) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_n, Y_1) & \text{cov}(X_n, Y_2) & \dots & \text{cov}(X_n, Y_m) \end{pmatrix}.$$

- We can also write:

$$\begin{aligned} \text{cov}(\mathbf{X}, \mathbf{Y}) &= E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{Y} - E[\mathbf{Y}])^T] \\ &= E \left[\begin{pmatrix} X_1 - E[X_1] \\ \vdots \\ X_n - E[X_n] \end{pmatrix} (Y_1 - E[Y_1], \dots, Y_m - E[Y_m]) \right]. \end{aligned}$$

- Some properties of covariance matrices between two vectors:
 - If \mathbf{A} and \mathbf{B} are constant matrices, $\text{cov}(\mathbf{AX}, \mathbf{BY}) = \mathbf{A} \text{cov}(\mathbf{X}, \mathbf{Y}) \mathbf{B}^T$.
 - If $\mathbf{Z} = \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}$, then $\text{cov}(\mathbf{Z}) = \begin{pmatrix} \text{cov}(\mathbf{X}) & \text{cov}(\mathbf{X}, \mathbf{Y}) \\ \text{cov}(\mathbf{Y}, \mathbf{X}) & \text{cov}(\mathbf{Y}) \end{pmatrix}$.
- Let \mathbf{X} be a random vector. The *correlation matrix* of \mathbf{X} is defined as:

$$\text{corr}(\mathbf{X}) = (\text{corr}(X_i, X_j)) = \begin{pmatrix} 1 & \text{corr}(X_1, X_2) & \dots & \text{corr}(X_1, X_n) \\ \text{corr}(X_2, X_1) & 1 & \dots & \text{corr}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{corr}(X_n, X_1) & \text{corr}(X_n, X_2) & \dots & 1 \end{pmatrix},$$

where

$$\text{corr}(X_i, X_j) = \frac{\text{cov}(X_i, X_j)}{\sqrt{\text{var}(X_i)\text{var}(X_j)}}.$$

1.3 Generalised Linear Models: A Review

Generalised linear models (GLMs) are developed by Nelder and Wedderburn [1972] as a way to unify other statistical models as well as the statistical methods operated on these models. In this section, we will briefly review the definition of GLMs. Please refer to the last term's lecture notes for more detailed derivations and examples.

Definition. A GLM is specified through the following components:

- A *linear predictor*: $\eta = \boldsymbol{\beta}^T \mathbf{x}$.
- An *injective response function* h , such that $\mu = E[Y|\mathbf{x}, \boldsymbol{\beta}] = h(\eta) = h(\boldsymbol{\beta}^T \mathbf{x})$. Equivalently, one can write $g(\mu) = \boldsymbol{\beta}^T \mathbf{x}$, where $g = h^{-1}$ is the *link* function.
- The *distributional assumption*: $P(Y|\mathbf{x}, \boldsymbol{\beta})$ is an EDF, that is:

$$P(y|\mathbf{x}, \boldsymbol{\beta}) = P(y|\theta(\mathbf{x}, \boldsymbol{\beta}), \phi(\mathbf{x}, \boldsymbol{\beta})) = \exp \left(\frac{y\theta - b(\theta)}{\phi} + c(y, \phi) \right). \quad (1.1)$$

From the properties of the EDF, the mean and variance of this distribution are:

$$E[Y|\theta, \phi] = \mu = b'(\theta) \quad (1.2)$$

$$\text{Var}[Y|\theta, \phi] = \phi b''(\theta) = \phi b''((b')^{-1}(\mu)) = \phi \mathcal{V}(\mu). \quad (1.3)$$

- We also assume *independent data*, that is:

$$P(\{y_i\} | \{\mathbf{x}_i\}, \boldsymbol{\beta}) = \prod_{i=1}^n P(y_i | \mathbf{x}_i, \boldsymbol{\beta}) \quad (1.4)$$

where $\{y_i, i = 1, \dots, n\}$ are response data given the $\{\mathbf{x}_i, i = 1, \dots, n\}$.

The Natural/Canonical Link. Recall that we have both:

$$\mu = E[Y|\theta, \phi] = b'(\theta) \quad (1.5)$$

$$\mu = E[Y|\mathbf{x}, \boldsymbol{\beta}] = h(\boldsymbol{\beta}^T \mathbf{x}) = h(\eta) \quad (1.6)$$

with Equation (1.5) holding as a result of $P(y|\theta, \phi)$ following an EDF distribution, and Equation (1.6) holding by definition for a GLM.

The *natural link* is the choice $h = b'$, or equivalently $g = (b')^{-1}$, resulting in the equation

$$\theta = \boldsymbol{\beta}^T \mathbf{x} = \eta. \quad (1.7)$$

1.4 Exercises

1. Prove the identities in Section 1.2.
2. The table below gives some common link functions for GLMs. In the table, $\Phi(\cdot)$ is the cdf of the standard normal distribution. Find their inverses (i.e., the response function h).

Link function	$\eta_i = g(\mu_i)$
Identity	μ_i
Log	$\log(\mu_i)$
Inverse	μ_i^{-1}
Inverse-square	μ_i^{-2}
Square-root	$\sqrt{\mu_i}$
Logit	$\log \frac{\mu_i}{1-\mu_i}$
Probit	$\Phi^{-1}(\mu_i)$
Log-log	$-\log(-\log \mu_i)$
Complementary log-log	$\log(-\log(1 - \mu_i))$

3. Write down the GLM components and the corresponding natural link when the response variable Y is assumed to follow each distribution below. Note that the values of y_i have to be in the correct range for each distribution.
 - Gaussian: $y_i | \mathbf{x}_i, \boldsymbol{\beta} \sim \mathcal{N}(\mu_i, \sigma^2)$.

- Bernoulli: $y_i|\mathbf{x}_i, \boldsymbol{\beta} \sim \text{Bernoulli}(\mu_i)$.
- Binomial: $y_i|\mathbf{x}_i, \boldsymbol{\beta} \sim \text{Bin}(n_i, \mu_i)$. Note that here y_i is the proportion of successes in n_i independent trials with probability of success μ_i . Thus, the range of y_i is $\{\frac{0}{n_i}, \frac{1}{n_i}, \dots, \frac{n_i}{n_i}\}$.
- Poisson: $y_i|\mathbf{x}_i, \boldsymbol{\beta} \sim \text{Poi}(\mu_i)$.
- Gamma: $y_i|\mathbf{x}_i, \boldsymbol{\beta} \sim \text{Gamma}(\nu, \alpha_i)$, where ν and α_i are the shape and rate parameters of the Gamma distribution.
- Inverse-Gaussian: $y_i|\mathbf{x}_i, \boldsymbol{\beta} \sim \text{IG}(\mu_i, \lambda)$.

From these results, show that the linear regression model, logistic regression model, and Poisson regression model are all special cases of GLMs.

4. In \mathbf{R} , find the default link function for each GLM family. Are they the natural link for the corresponding family?

Chapter 2

Estimation

Suppose we are given a dataset that we would like to model using a GLM. After checking the data and possibly performing some exploratory data analysis, we have already chosen a specific form of GLMs that is most suitable for the dataset (say e.g., the Poisson GLM). Now we wish to estimate the parameters of the model; that is, finding the value of β that best explains our data. For this purpose, we can use the maximum likelihood method to obtain an estimate $\hat{\beta}$.

One advantage of the GLM formulation is that it allows us to derive a solution for the maximum likelihood method directly in the general form, thus unifying the estimation methods for various statistical models. In other words, we do not need to derive solutions for the Poisson GLM, the Binomial GLM, etc. separately. Instead, we can just derive the solution for the general form of the GLMs. The solution for each specific form of the GLM will be an instance of this general solution.

For the rest of this chapter,

2.1 Likelihood Function

Consider the grouped data setup where we have predictors and data with possible replicates $\{(\mathbf{x}_i, y_{ir_i})\}_{i \in [1..n], r_i \in [1..m_i]}$. Recall that under a GLM, given predictors $\{\mathbf{x}_i\}_{i \in [1..n]}$, each response y_{ir_i} is independent of the other y_{jr_j} , and of the values of all predictors \mathbf{x}_j with $j \neq i$, so that the joint probability of the data — that is, the likelihood — is given by

$$L(\beta) = P(\{y_{ir_i}\} | \{\mathbf{x}_i\}, \beta) = P(\{y_{ir_i}\} | \{\theta_i\}, \phi) = \prod_{i=1}^n \prod_{r_i=1}^{m_i} P(y_{ir_i} | \theta_i, \phi) \quad (2.1)$$

where

$$P(y_{ir_i} | \theta_i, \phi) = \exp \left(\frac{y_{ir_i} \theta_i - b(\theta_i)}{\phi} + c(y_{ir_i}, \phi) \right) \quad (2.2)$$

with

$$\theta_i = (b')^{-1}(\mu_i) = (b')^{-1}(h(\eta_i)) = (b')^{-1}(h(\beta^T \mathbf{x}_i)). \quad (2.3)$$

2.2 Log-Likelihood Function

The log probability of the data — or log-likelihood — is thus given by

$$l(\beta) = \log L(\beta) = \log P(\{y_{ir_i}\} | \{\theta_i\}, \phi) \quad (2.4)$$

$$= \sum_i \sum_{r_i} \left(\frac{y_{ir_i} \theta_i - b(\theta_i)}{\phi} + c(y_{ir_i}, \phi) \right) \quad (2.5)$$

$$= \sum_i \left(m_i \frac{y_i \theta_i - b(\theta_i)}{\phi} + \sum_{r_i} c(y_{ir_i}, \phi) \right) \quad (2.6)$$

$$= \sum_i l_i \quad (2.7)$$

where we have defined

$$y_i = \frac{1}{m_i} \sum_{r_i} y_{ir_i} \quad (2.8)$$

$$l_i = \frac{y_i \theta_i - b(\theta_i)}{\phi_i} + \sum_{r_i} c(y_{ir_i}, \phi) \quad (2.9)$$

$$\phi_i = \phi / m_i. \quad (2.10)$$

2.3 Score Function and Score Equation

The *score function* is given by

$$\mathbf{S}(\beta) = \frac{\partial l}{\partial \beta^T} = \sum_i \frac{\partial l_i}{\partial \beta^T} = \sum_i \frac{\partial l_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta^T} \quad (2.11)$$

where, recalling that $\mu_i = b'(\theta_i)$, $\mathcal{V}(\mu_i) = b''(\theta_i)$, $\mu_i = h(\eta_i)$ and $\eta_i = \beta^T \mathbf{x}_i$, we have:¹

$$\frac{\partial l_i}{\partial \theta_i} = \frac{y_i - b'(\theta_i)}{\phi_i} = \frac{y_i - \mu_i}{\phi_i} \quad (2.12)$$

$$\frac{\partial \theta_i}{\partial \mu_i} = 1 / \left(\frac{\partial \mu_i}{\partial \theta_i} \right) = \frac{1}{b''(\theta_i)} = \frac{1}{\mathcal{V}(\mu_i)} \quad (2.13)$$

$$\frac{\partial \mu_i}{\partial \eta_i} = h'(\eta_i) \quad (2.14)$$

$$\frac{\partial \eta_i}{\partial \beta^T} = \mathbf{x}_i. \quad (2.15)$$

¹Note that here and in this module, we assume that ϕ does not depend on β .

The score function is thus given by

$$\mathbf{S}(\boldsymbol{\beta}) = \sum_i \left(\frac{y_i - \mu_i}{\phi_i} \right) \left(\frac{1}{\mathcal{V}(\mu_i)} \right) h'(\eta_i) \mathbf{x}_i \quad (2.16)$$

$$= \frac{1}{\phi} \sum_i m_i(y_i - \mu_i) \frac{1}{\mathcal{V}(\mu_i)} h'(\eta_i) \mathbf{x}_i. \quad (2.17)$$

The maximum likelihood estimate $\hat{\boldsymbol{\beta}}$ must then satisfy the *score equation*:

$$\mathbf{S}(\hat{\boldsymbol{\beta}}) = 0. \quad (2.18)$$

Note that the dispersion parameter ϕ cancels from the score equation, which implies that $\hat{\boldsymbol{\beta}}$ does not depend on ϕ . This is another important property of EDFs.

2.3.1 Special Case: Natural Link

For the natural link, $\theta_i = \eta_i$, so Equations (2.13) and (2.14) combine to give

$$\frac{h'(\eta_i)}{\mathcal{V}(\mu_i)} = \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} = \frac{\partial \theta_i}{\partial \eta_i} = 1. \quad (2.19)$$

The score function thus simplifies to

$$\mathbf{S}(\boldsymbol{\beta}) = \frac{1}{\phi} \sum_i m_i(y_i - \mu_i) \mathbf{x}_i. \quad (2.20)$$

2.4 Fisher Information

To solve the score equation, we will also need the second derivative of the log-likelihood. Its negative is called the *Observed Fisher Information*, defined as

$$\mathbf{F}_{\text{obs}}(\boldsymbol{\beta}) = -\frac{\partial^2 l}{\partial \boldsymbol{\beta}^T \partial \boldsymbol{\beta}} = -\frac{\partial \mathbf{S}}{\partial \boldsymbol{\beta}}. \quad (2.21)$$

Note that, at the MLE, $\mathbf{F}_{\text{obs}}(\hat{\boldsymbol{\beta}})$ is positive by definition. Because it is a function of the data $\{y_i\}$, \mathbf{F}_{obs} has a probability distribution. In practice, the Observed Fisher Information is often approximated by the *Expected Fisher Information*, otherwise known simply as the *Fisher Information*.²

$$\mathbf{F}(\boldsymbol{\beta}) = E \left[-\frac{\partial \mathbf{S}}{\partial \boldsymbol{\beta}} \right] \quad (2.22)$$

²Some texts refer to $\mathbf{F}_{\text{obs}}(\hat{\boldsymbol{\beta}})$ as the Observed Fisher Information, and to $\mathbf{F}_{\text{obs}}(\boldsymbol{\beta})$ simply as the Fisher Information. Some don't refer to either of these at all. Just to be clear, we will refer to $\mathbf{F}_{\text{obs}}(\boldsymbol{\beta})$ as the Observed Fisher Information and the Expected Fisher Information $\mathbf{F}(\boldsymbol{\beta})$ simply as the Fisher Information.

where the expectation is taken over the joint probability distribution of the data $P(\{y_{ir_i}\}|\boldsymbol{\beta}, \{\mathbf{x}_i\})$.

2.5 Example: Poisson Regression

We look at two example calculations of the score function and Fisher Information for Poisson Regression, that is we have

- $y|\mathbf{x}, \boldsymbol{\beta} \sim \text{Poi}(\lambda(\mathbf{x}, \boldsymbol{\beta}))$
- $\phi = 1$.

In this example, let us assume $m_i = 1$, i.e., the data are ungrouped.

2.5.1 With Natural Link

We have that: $\lambda(\mathbf{x}, \boldsymbol{\beta}) = \mu(\mathbf{x}, \boldsymbol{\beta}) = h(\eta(\mathbf{x}, \boldsymbol{\beta})) = e^{\eta(\mathbf{x}, \boldsymbol{\beta})} = e^{\boldsymbol{\beta}^T \mathbf{x}}$.

Equation (2.20) then gives

$$\mathbf{S}(\boldsymbol{\beta}) = \sum_i (y_i - e^{\boldsymbol{\beta}^T \mathbf{x}_i}) \mathbf{x}_i \quad (2.23)$$

while Equation (2.21) gives

$$\mathbf{F}_{\text{obs}}(\boldsymbol{\beta}) = \sum_i e^{\boldsymbol{\beta}^T \mathbf{x}_i} \mathbf{x}_i \mathbf{x}_i^T. \quad (2.24)$$

Note that this does not depend on the data, so that Equation (2.22) gives

$$\mathbf{F}(\boldsymbol{\beta}) = \mathbb{E}[\mathbf{F}_{\text{obs}}(\boldsymbol{\beta})] = \mathbf{F}_{\text{obs}}(\boldsymbol{\beta}). \quad (2.25)$$

2.5.2 With Identity Link

The identity link is defined such that $h(\eta) = \eta$.

In this case, we have that:

- $\lambda(\mathbf{x}, \boldsymbol{\beta}) = \mu(\mathbf{x}, \boldsymbol{\beta}) = h(\eta(\mathbf{x}, \boldsymbol{\beta})) = \eta(\mathbf{x}, \boldsymbol{\beta}) = \boldsymbol{\beta}^T \mathbf{x}$.
- $\mathcal{V}(\mu) = \mu$ (see Poisson example in EDF chapter).
- $h'(\eta) = 1$.

Equation (2.17) thus gives

$$\mathbf{S}(\boldsymbol{\beta}) = \sum_i (y_i - \mu_i) \frac{1}{\mu_i} \mathbf{1} \mathbf{x}_i \quad (2.26)$$

$$= \sum_i (y_i - \boldsymbol{\beta}^T \mathbf{x}_i) \frac{1}{\boldsymbol{\beta}^T \mathbf{x}_i} \mathbf{x}_i \quad (2.27)$$

$$= \sum_i \left(\frac{y_i}{\boldsymbol{\beta}^T \mathbf{x}_i} - 1 \right) \mathbf{x}_i. \quad (2.28)$$

Using Equation (2.21) and the chain rule, we can obtain the Observed Fisher Information:

$$\mathbf{F}_{\text{obs}}(\boldsymbol{\beta}) = \sum_i \frac{y_i}{(\boldsymbol{\beta}^T \mathbf{x}_i)^2} \mathbf{x}_i \mathbf{x}_i^T. \quad (2.29)$$

Hence, the Fisher Information is:

$$\mathbf{F}(\boldsymbol{\beta}) = \mathbb{E}[\mathbf{F}_{\text{obs}}(\boldsymbol{\beta})] \quad (2.30)$$

$$= \mathbb{E} \left[\sum_i \frac{Y_i}{(\boldsymbol{\beta}^T \mathbf{x}_i)^2} \mathbf{x}_i \mathbf{x}_i^T \right] \quad (2.31)$$

$$= \sum_i \frac{\mathbb{E}[Y_i | \boldsymbol{\beta}, \mathbf{x}_i]}{(\boldsymbol{\beta}^T \mathbf{x}_i)^2} \mathbf{x}_i \mathbf{x}_i^T \quad (2.32)$$

$$= \sum_i \frac{\boldsymbol{\beta}^T \mathbf{x}_i}{(\boldsymbol{\beta}^T \mathbf{x}_i)^2} \mathbf{x}_i \mathbf{x}_i^T \quad (2.33)$$

$$= \sum_i \frac{1}{\boldsymbol{\beta}^T \mathbf{x}_i} \mathbf{x}_i \mathbf{x}_i^T. \quad (2.34)$$

Note that $\mathbf{F}(\boldsymbol{\beta}) \neq \mathbf{F}_{\text{obs}}(\boldsymbol{\beta})$ in this case.

2.6 Properties of $\mathbf{S}(\boldsymbol{\beta})$ and $\mathbf{F}(\boldsymbol{\beta})$

Having defined the score function $\mathbf{S}(\boldsymbol{\beta})$ and the Fisher information $\mathbf{F}(\boldsymbol{\beta})$, in this section we will investigate some of their properties. First, let us define $S_i(\boldsymbol{\beta}) = \frac{\partial l_i}{\partial \boldsymbol{\beta}}$. Then we have $\mathbf{S}(\boldsymbol{\beta}) = \sum_i \mathbf{S}_i(\boldsymbol{\beta})$.

2.6.1 Expectation of $\mathbf{S}(\boldsymbol{\beta})$

The expectation of $\mathbf{S}(\boldsymbol{\beta})$ can be computed from Equations (2.11) and (2.17) as follows:

$$\mathbb{E}[\mathbf{S}(\boldsymbol{\beta})] = \sum_i \mathbb{E}[\mathbf{S}_i(\boldsymbol{\beta})] \quad (2.35)$$

$$= \sum_i \frac{\mathbb{E}[Y_i | \boldsymbol{\beta}, \mathbf{x}_i] - \mu_i}{\phi_i} \frac{1}{\mathcal{V}(\mu_i)} h'(\eta_i) \mathbf{x}_i \quad (2.36)$$

$$= 0 \quad (2.37)$$

because $E[Y_i|\boldsymbol{\beta}, \mathbf{x}_i] = \mu_i$.

2.6.2 Variance of $\mathbf{S}(\boldsymbol{\beta})$

Using Equation (2.17) and the properties of covariance matrices in Section 1.2, we can calculate the variance of $\mathbf{S}(\boldsymbol{\beta})$ as follows:

$$\text{Var}[\mathbf{S}(\boldsymbol{\beta})] = \sum_i \text{Var}[\mathbf{S}_i(\boldsymbol{\beta})] \quad (2.38)$$

$$= \sum_i \text{Var} \left[\frac{h'(\eta_i)}{\phi_i \mathcal{V}(\mu_i)} \mathbf{x}_i (Y_i - \mu_i) \right] \quad (2.39)$$

$$= \sum_i \left(\frac{h'(\eta_i)}{\phi_i \mathcal{V}(\mu_i)} \mathbf{x}_i \right) \text{Var}[Y_i - \mu_i] \left(\frac{h'(\eta_i)}{\phi_i \mathcal{V}(\mu_i)} \mathbf{x}_i^T \right) \quad (2.40)$$

$$= \sum_i \left(\frac{h'(\eta_i)^2}{\phi_i^2 \mathcal{V}(\mu_i)^2} \mathbf{x}_i \mathbf{x}_i^T \right) \text{Var}[Y_i]. \quad (2.41)$$

Note that in the first equality above, we can break the variance into sum of smaller components due to the independent data assumption.

Now using $\text{Var}[Y_i] = \phi_i \mathcal{V}(\mu_i)$, we can obtain the expression for $\text{Var}[\mathbf{S}(\boldsymbol{\beta})]$:

$$\text{Var}[\mathbf{S}(\boldsymbol{\beta})] = \sum_i \frac{h'(\eta_i)^2}{\phi_i \mathcal{V}(\mu_i)} \mathbf{x}_i \mathbf{x}_i^T. \quad (2.42)$$

2.6.3 Properties of $\mathbf{F}(\boldsymbol{\beta})$

Recall from Equation (2.22) that $\mathbf{F}(\boldsymbol{\beta}) = E \left[-\frac{\partial \mathbf{S}}{\partial \boldsymbol{\beta}} \right] = -E \left[\frac{\partial^2 l}{\partial \boldsymbol{\beta}^T \partial \boldsymbol{\beta}} \right]$.

We will first show that $E \left[\frac{\partial^2 l}{\partial \boldsymbol{\beta}^T \partial \boldsymbol{\beta}} \right] = E \left[-\frac{\partial l}{\partial \boldsymbol{\beta}^T} \frac{\partial l}{\partial \boldsymbol{\beta}} \right]$.

2.6.3.1 An Important Identity

Let $\rho = e^l$, where l is the log-likelihood, so that $\rho = L(\boldsymbol{\beta}) = P(\{y_{ir_i}\} | \{\mathbf{x}_i\}, \boldsymbol{\beta})$ is the likelihood/probability of the data. Then

$$\frac{\partial l}{\partial \boldsymbol{\beta}^T} = \frac{\partial l}{\partial \rho} \frac{\partial \rho}{\partial \boldsymbol{\beta}^T} = \frac{1}{\rho} \frac{\partial \rho}{\partial \boldsymbol{\beta}^T}. \quad (2.43)$$

Using the product rule and chain rule, we have

$$\frac{\partial^2 l}{\partial \beta^T \partial \beta} = -\frac{1}{\rho^2} \frac{\partial \rho}{\partial \beta^T} \frac{\partial \rho}{\partial \beta} + \frac{1}{\rho} \frac{\partial^2 \rho}{\partial \beta^T \partial \beta} \quad (2.44)$$

$$= -\frac{1}{\rho^2} \left(\frac{\partial \rho}{\partial l} \frac{\partial l}{\partial \beta^T} \right) \left(\frac{\partial \rho}{\partial l} \frac{\partial l}{\partial \beta} \right) + \frac{1}{\rho} \frac{\partial^2 \rho}{\partial \beta^T \partial \beta} \quad (2.45)$$

$$= -\frac{1}{\rho^2} \left(\rho \frac{\partial l}{\partial \beta^T} \right) \left(\rho \frac{\partial l}{\partial \beta} \right) + \frac{1}{\rho} \frac{\partial^2 \rho}{\partial \beta^T \partial \beta} \quad (2.46)$$

$$= -\frac{\partial l}{\partial \beta^T} \frac{\partial l}{\partial \beta} + \frac{1}{\rho} \frac{\partial^2 \rho}{\partial \beta^T \partial \beta}. \quad (2.47)$$

Note that the expectation (over the data) of the second term is

$$\mathbb{E} \left[\frac{1}{\rho} \frac{\partial^2 \rho}{\partial \beta^T \partial \beta} \right] = \int \rho \frac{1}{\rho} \frac{\partial^2 \rho}{\partial \beta^T \partial \beta} = \int \frac{\partial^2 \rho}{\partial \beta^T \partial \beta} = \frac{\partial^2}{\partial \beta^T \partial \beta} \int \rho = \frac{\partial^2}{\partial \beta^T \partial \beta} 1 = 0. \quad (2.48)$$

Thus, we have:

$$E \left[\frac{\partial^2 l}{\partial \beta^T \partial \beta} \right] = E \left[-\frac{\partial l}{\partial \beta^T} \frac{\partial l}{\partial \beta} \right]. \quad (2.49)$$

2.6.3.2 Relating $\mathbf{F}(\beta)$ and $\text{Var}[\mathbf{S}(\beta)]$

Using Equation (2.49), we have that:

$$\mathbf{F}(\beta) = -\mathbb{E} \left[\frac{\partial^2 l}{\partial \beta^T \partial \beta} \right] = \mathbb{E} \left[\frac{\partial l}{\partial \beta^T} \frac{\partial l}{\partial \beta} \right] \quad (2.50)$$

$$= \mathbb{E}[\mathbf{S}(\beta) \mathbf{S}(\beta)^T] \quad (2.51)$$

$$= \text{Var}[\mathbf{S}(\beta)] + \mathbb{E}[\mathbf{S}(\beta)] \mathbb{E}[\mathbf{S}(\beta)]^T \quad (2.52)$$

$$= \text{Var}[\mathbf{S}(\beta)] \quad (2.53)$$

where the last equality is due to $\mathbb{E}[\mathbf{S}(\beta)] = 0$.

Therefore, an important property of the Fisher Information is that it is equal to the variance of the score function, whose expression is given in Equation (2.42).

2.6.3.3 Special Case: Natural Link

For the natural link, recall that $\frac{h'(\eta_i)}{\mathcal{V}(\mu_i)} = 1$. So we have:

$$\mathbf{S}(\beta) = \sum_i \frac{1}{\phi_i} (y_i - h(\eta_i)) \mathbf{x}_i. \quad (2.54)$$

Let $\mathbf{S}_i = \frac{1}{\phi_i} (y_i - h(\eta_i)) \mathbf{x}_i$. We have $\mathbf{S} = \sum_i \mathbf{S}_i$ and thus:

$$\mathbf{F}_{\text{obs}}(\boldsymbol{\beta}) = -\frac{\partial \mathbf{S}}{\partial \boldsymbol{\beta}} = -\sum_i \frac{\partial \mathbf{S}_i}{\partial \boldsymbol{\beta}} = -\sum_i \frac{\partial \mathbf{S}_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \boldsymbol{\beta}} = \sum_i \frac{h'(\eta_i)}{\phi_i} \mathbf{x}_i \mathbf{x}_i^T \quad (2.55)$$

$$\mathbf{F}(\boldsymbol{\beta}) = \text{Var}[\mathbf{S}(\boldsymbol{\beta})] = \sum_i \frac{h'(\eta_i)}{\phi_i} \mathbf{x}_i \mathbf{x}_i^T. \quad (2.56)$$

Thus, for the natural link, we see that $\mathbf{F}(\boldsymbol{\beta}) = \mathbf{F}_{\text{obs}}(\boldsymbol{\beta})$.

2.7 Matrix Notation

For the next section, it is useful to establish a condensed, matrix notation for some of the previous quantities, analogous to the matrix notation used for linear models.

- Let $\mathbf{Y} \in \mathbb{R}^n$ be the random vector with components Y_i , the response values. This is exactly the same quantity as in the linear model case.
- Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ be the *design matrix*, the matrix with components $x_{i,a}$, the value of the a^{th} component of the predictor vector for the i^{th} data point. This is exactly the same quantity as in the linear model case. This matrix is sometimes called the *model matrix*.
- Let $\boldsymbol{\mu} \in \mathbb{R}^n$ be the vector with components $\mu_i = h(\boldsymbol{\beta}^T \mathbf{x}_i)$, so that $\boldsymbol{\mu} = \mathbb{E}[\mathbf{Y}]$.
- Let $\mathbf{D} \in \mathbb{R}^{n \times n}$ be the diagonal matrix with components $D_{ii} = h'(\eta_i)$. For example, if $h(\eta) = e^\eta$, then

$$\mathbf{D} = \begin{pmatrix} e^{\boldsymbol{\beta}^T \mathbf{x}_1} & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & e^{\boldsymbol{\beta}^T \mathbf{x}_n} \end{pmatrix}. \quad (2.57)$$

- Let $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$ be the covariance matrix for \mathbf{Y} , with components:

$$\Sigma_{ij} = \text{Cov}[Y_i, Y_j] = \text{Var}[Y_i] \delta_{ij} = \phi_i \mathcal{V}(\mu_i) \delta_{ij}. \quad (2.58)$$

That is,

$$\boldsymbol{\Sigma} = \begin{pmatrix} \text{Var}[Y_1] & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \text{Var}[Y_n] \end{pmatrix} = \begin{pmatrix} \phi_1 \mathcal{V}(\mu_1) & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & \phi_n \mathcal{V}(\mu_n) \end{pmatrix}. \quad (2.59)$$

2.7.1 Score Function and Fisher Information

Recall that

$$\mathbf{S}(\boldsymbol{\beta}) = \sum_i \left(\frac{y_i - \mu_i}{\phi_i \mathcal{V}(\mu_i)} \right) h'(\eta_i) \mathbf{x}_i \quad (2.60)$$

$$\mathbf{F}(\boldsymbol{\beta}) = \sum_i \frac{h'(\eta_i)^2}{\phi_i \mathcal{V}(\mu_i)} \mathbf{x}_i \mathbf{x}_i^T. \quad (2.61)$$

In terms of the matrix notation, these become

$$\mathbf{S} = \mathbf{X}^T \mathbf{D} \boldsymbol{\Sigma}^{-1} (\mathbf{Y} - \boldsymbol{\mu}) \quad (2.62)$$

$$\mathbf{F} = \mathbf{X}^T \mathbf{D}^T \boldsymbol{\Sigma}^{-1} \mathbf{D} \mathbf{X}. \quad (2.63)$$

2.7.2 Natural Link

Note that for the natural link,

$$\frac{\partial \theta_i}{\partial \eta_i} = \frac{h'(\eta_i)}{\mathcal{V}(\mu_i)} = 1. \quad (2.64)$$

Thus, with $\phi_i = \phi/m_i$, we have:

$$h'(\eta_i) = \mathcal{V}(\mu_i) = \frac{\text{Var}[Y_i]}{\phi_i} = m_i \frac{\text{Var}[Y_i]}{\phi}. \quad (2.65)$$

Now let $\mathbf{G} \in \mathbb{R}^{n \times n}$ be the diagonal matrix with components $m_i \delta_{ij}$, known as the *grouping* matrix. Then

$$\mathbf{D} = \frac{1}{\phi} \mathbf{G} \boldsymbol{\Sigma} = \frac{1}{\phi} \boldsymbol{\Sigma} \mathbf{G}, \quad (2.66)$$

and therefore,

$$\mathbf{S}(\boldsymbol{\beta}) = \frac{1}{\phi} \mathbf{X}^T \mathbf{G} (\mathbf{Y} - \boldsymbol{\mu}) \quad (2.67)$$

$$\mathbf{F}(\boldsymbol{\beta}) = \frac{1}{\phi^2} \mathbf{X}^T \mathbf{G}^T \boldsymbol{\Sigma} \mathbf{G} \mathbf{X}. \quad (2.68)$$

2.8 Iterative Solution of $\mathbf{S}(\hat{\boldsymbol{\beta}}) = 0$

So far we have seen how to set up the score equation for the maximum likelihood estimate. We have also proven some of its properties as well as those of the Fisher Information. We now turn to the question of how to solve the score equation. In general, this cannot be done in closed form, except in rare cases. So we need to turn to numerical methods implemented on a computer.

We have the same two options here as in the binary regression case: we can try to optimise l directly, or we can attempt to solve the score equation. There are many algorithms that can be used to perform these tasks. Here we focus on one: *Iteratively Reweighted Least Squares (IRLS)*, also known as *Iterative Weighted Least Squares (IWLS)*.³

We start by recalling the Newton-Raphson method for finding the zero of a function. Note that we wish to solve the equation:

$$\mathbf{S}(\hat{\beta}) = 0. \quad (2.69)$$

We then approximate \mathbf{S} linearly about some point β_0 :

$$\mathbf{S}(\beta_0 + \delta\beta_0) = \mathbf{S}(\beta_0) + \frac{\partial \mathbf{S}(\beta_0)}{\partial \beta} \delta\beta_0 + \mathcal{O}(\delta\beta_0^2) \quad (2.70)$$

where the reason for the subscript 0 will become apparent soon. In the case when $\mathbf{S}(\beta_0 + \delta\beta_0) = 0$ (such as the case we are interested in), we have approximately that:

$$\frac{\partial \mathbf{S}(\beta_0)}{\partial \beta} \delta\beta_0 = -\mathbf{S}(\beta_0). \quad (2.71)$$

Now in our case,

$$-\frac{\partial \mathbf{S}(\beta_0)}{\partial \beta} = \mathbf{F}_{\text{obs}}(\beta_0), \quad (2.72)$$

so Equation (2.71) becomes

$$\mathbf{F}_{\text{obs}}(\beta_0) \delta\beta_0 = \mathbf{S}(\beta_0), \quad (2.73)$$

or equivalently,

$$\delta\beta_0 = (\mathbf{F}_{\text{obs}}(\beta_0))^{-1} \mathbf{S}(\beta_0). \quad (2.74)$$

This then gives us a new value:

$$\beta_1 = \beta_0 + \delta\beta_0 = \beta_0 + (\mathbf{F}_{\text{obs}}(\beta_0))^{-1} \mathbf{S}(\beta_0). \quad (2.75)$$

Then we can iterate the steps above for $m = 1, 2, \dots$,

$$\beta_{m+1} = \beta_m + \delta\beta_m \quad (2.76)$$

³You have cause to be particularly interested in this algorithm as a Durham student. It is on the undergraduate syllabus of nearly every maths degree in the world which includes a large statistical component and some of the important early development was researched by Dr Peter Green when he was a lecturer at Durham: Green [1984].

where

$$\delta\beta_m = (\mathbf{F}_{\text{obs}}(\beta_m))^{-1} \mathbf{S}(\beta_m). \quad (2.77)$$

Because \mathbf{F}_{obs} is hard to find and hard to invert in general, we approximate it with the expected Fisher Information. This is known as the *Fisher scoring* method, where we compute $\delta\beta_m$ by:

$$\delta\beta_m = (\mathbf{F}(\beta_m))^{-1} \mathbf{S}(\beta_m). \quad (2.78)$$

2.8.1 Iteratively Reweighted Least Squares (IRLS)

Now we will use Equation (2.78) to derive the IRLS method in matrix notation. From Equation (2.78), we have that

$$\mathbf{F}(\beta_m)\delta\beta_m = \mathbf{S}(\beta_m) \quad (2.79)$$

or equivalently that

$$\mathbf{F}(\beta_m)\beta_{m+1} = \mathbf{F}(\beta_m)\beta_m + \mathbf{S}(\beta_m). \quad (2.80)$$

Using the matrix notation in Section 2.7.1 and defining $\mathbf{W} = \mathbf{D}^T \Sigma^{-1} \mathbf{D}$, we can write

$$\mathbf{F} = \mathbf{X}^T \mathbf{D}^T \Sigma^{-1} \mathbf{D} \mathbf{X} = \mathbf{X}^T \mathbf{W} \mathbf{X} \quad (2.81)$$

and

$$\mathbf{S} = \mathbf{X}^T \mathbf{D} \Sigma^{-1} (\mathbf{Y} - \boldsymbol{\mu}) = \mathbf{X}^T \mathbf{D}^T \Sigma^{-1} (\mathbf{Y} - \boldsymbol{\mu}) = \mathbf{X}^T \mathbf{W} \mathbf{D}^{-1} (\mathbf{Y} - \boldsymbol{\mu}), \quad (2.82)$$

since \mathbf{D} is a diagonal matrix.

Thus, using the subscript m to denote the value of a quantity evaluated using β_m or derived quantities, we can calculate the right hand side of Equation (2.80) from:

$$\mathbf{F}_m \beta_m + \mathbf{S}_m = \mathbf{X}^T \mathbf{W}_m \mathbf{X} \beta_m + \mathbf{X}^T \mathbf{W}_m \mathbf{D}_m^{-1} (\mathbf{Y} - \boldsymbol{\mu}_m) = \mathbf{X}^T \mathbf{W}_m \tilde{\mathbf{Y}}_m \quad (2.83)$$

where

$$\tilde{\mathbf{Y}}_m = \mathbf{X} \beta_m + \mathbf{D}_m^{-1} (\mathbf{Y} - \boldsymbol{\mu}_m) \quad (2.84)$$

are the so-called *working observations*.

Now replacing $\mathbf{F}(\beta_m) = \mathbf{F}_m = \mathbf{X}^T \mathbf{W}_m \mathbf{X}$ in the left hand side of Equation (2.80), we have that

$$(\mathbf{X}^T \mathbf{W}_m \mathbf{X}) \beta_{m+1} = \mathbf{X}^T \mathbf{W}_m \tilde{\mathbf{Y}}_m \quad (2.85)$$

or

$$\beta_{m+1} = (\mathbf{X}^T \mathbf{W}_m \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}_m \tilde{\mathbf{Y}}_m. \quad (2.86)$$

Thus, to find a solution for $\mathbf{S}(\beta) = 0$, we can start from an arbitrary point β_0 and iteratively apply Equation (2.86) until a convergence criterion is met.

This sequence of iterated operations is called *iteratively reweighted least squares* or *iterative weighted least squares* since each iteration is the solution to the following least squares problem: minimize the quantity $l_m(\beta)$ with respect to β , where

$$l_m(\beta) = (\tilde{\mathbf{Y}}_m - \mathbf{X}\beta)^T \mathbf{W}_m (\tilde{\mathbf{Y}}_m - \mathbf{X}\beta) \quad (2.87)$$

and \mathbf{W}_m is known as the *weight matrix*.

2.8.2 IRLS Pseudo-Code

We now give a pseudo-code below for running IRLS. Note that the code will not run without computing μ , \mathbf{D} and \mathbf{W} using a specific example.

```
IRLS <- function(Y, X, phi, epsilon) {
  # Pick an initial value for hatBeta
  hatbeta = initializeBeta()

  # Set up convergence
  converged = false

  # Loop as long as convergence condition is not satisfied
  while not converged loop
  {
    # Compute mu, D, and Sigma (use h, h', V as subroutines)
    mu = computeMu(hatBeta, X)
    D = computeD(hatBeta, X)
    Sigma = computeSigma(hatBeta, phi)

    # Compute the weight matrix W
    W = t(D) %*% solve(Sigma) %*% D

    # Compute the working observations tildeY
    tildeY = X %*% hatBeta + solve(D) %*% (Y - mu)

    # Compute the new value of hatBeta
    newHatBeta = solve(t(X) %*% W %*% X) %*% (t(X) %*% W %*% tildeY)
```



```

    # Check whether we have converged
    converged = ((norm(newHatBeta - hatBeta) / norm(hatBeta)) <= epsilon)

    # Store new value of hatBeta ready for next iteration or return
    hatBeta = newHatBeta
  }

  return hatBeta
}

```

2.9 Practical Example: US Polio Data

In this example, we will use the US Polio data discussed last term to fit the GLMs. Recall that this dataset is a matrix of count data, giving the monthly number of polio cases in the United States from 1970 to 1983. We will convert this dataset into a matrix with two columns:

- covariate `time` in the first column ranging from 1 to 168, starting with January 1970.
- response `cases` in the second column, indicating the monthly number of polio cases.

We now load the data from the library `gamlss.data` and do the conversion.

```

library("gamlss.data")

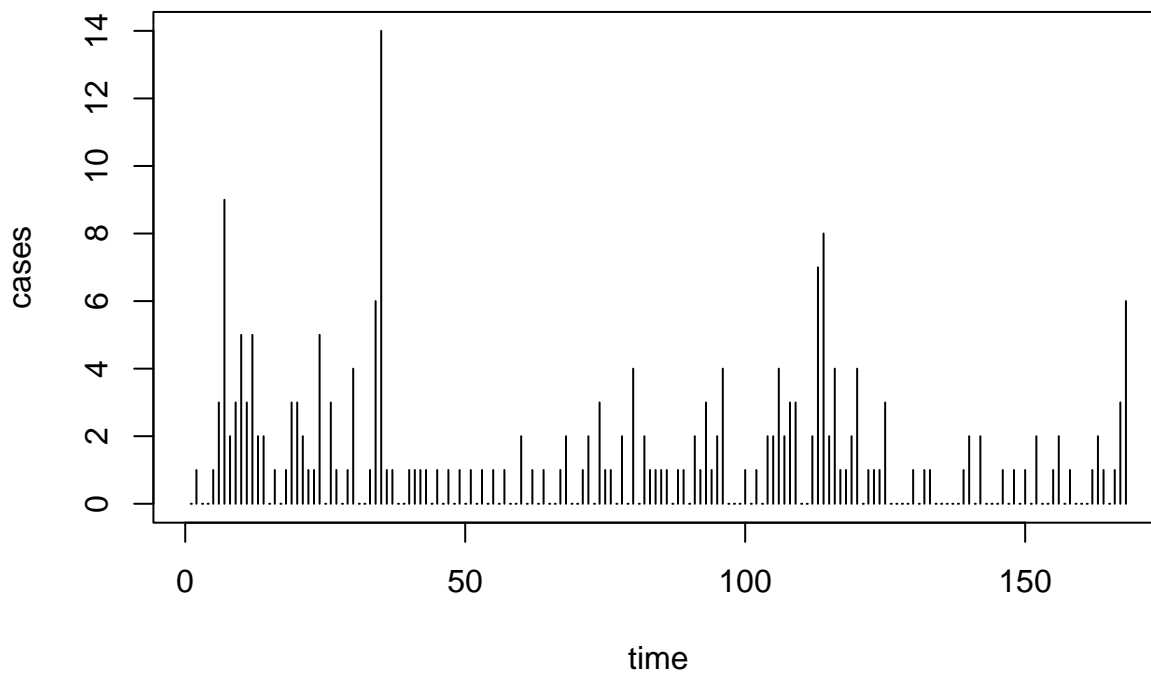
##
## Attaching package: 'gamlss.data'

## The following object is masked from 'package:datasets':
##
##      sleep
data("polio")
uspolio <- as.data.frame(matrix(c(1:168, t(polio)), ncol=2))
colnames(uspolio) <- c("time", "cases")

```

First, let us plot the data.

```
plot(uspolio, type="h")
```



Note that the main question we wish to consider is: *How has Polio incidence changed over time?*

Since this is count data, we begin by fitting a Poisson model with a linear time trend.

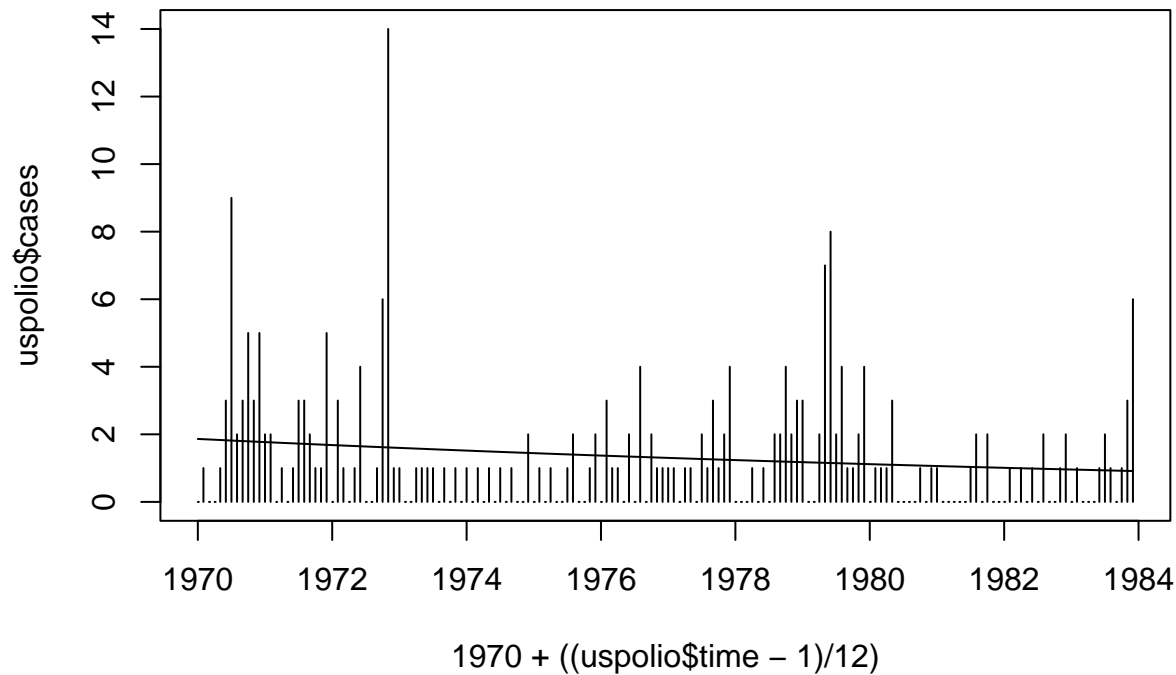
```
# Poisson model with linear time trend
polio.glm <- glm(cases ~ time, family=poisson(link=log), data=uspolio)

# Look at the model summary
summary(polio.glm)
```

```
##
## Call:
## glm(formula = cases ~ time, family = poisson(link = log), data = uspolio)
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.626639   0.123641   5.068 4.02e-07 ***
## time        -0.004263   0.001395  -3.055 0.00225 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
## Null deviance: 343.00  on 167  degrees of freedom
## Residual deviance: 333.55  on 166  degrees of freedom
## AIC: 594.59
##
## Number of Fisher Scoring iterations: 5
```

We can then plot the model as follows.

```
plot(1970 + ((uspolio$time - 1)/12), uspolio$cases, type="h")
lines(1970 + ((uspolio$time - 1)/12), polio.glm$fitted)
```



We can see that this is perhaps unsatisfactory. To improve the model, we can explore a linear trend with seasonal (annual) components.

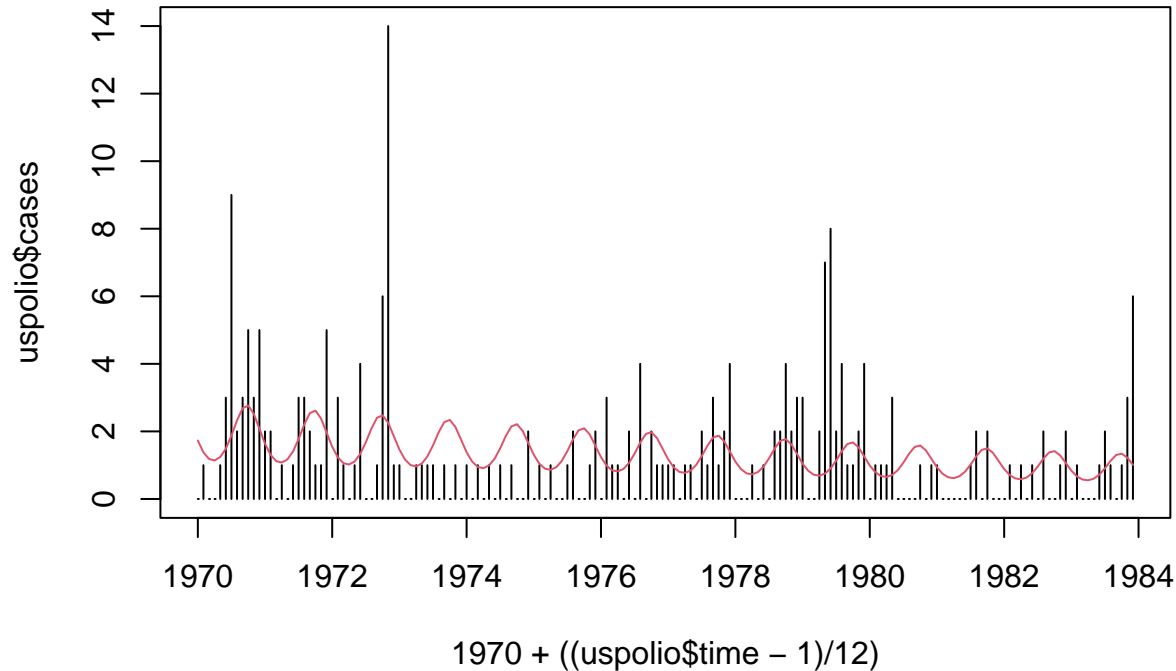
```
# Poisson model with linear trend and annual components
polio1.glm <- glm(cases ~ time + I(cos(2*pi*time/12)) + I(sin(2*pi*time/12)),
family=poisson(link=log), data=uspolio)
```

```
summary(polio1.glm)
```

```
##
## Call:
## glm(formula = cases ~ time + I(cos(2 * pi * time/12)) + I(sin(2 *
##   pi * time/12)), family = poisson(link = log), data = uspolio)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.606612   0.124800   4.861 1.17e-06 ***
## time          -0.004644   0.001401  -3.315 0.000916 ***
## I(cos(2 * pi * time/12))  0.181254   0.096160   1.885 0.059442 .
## I(sin(2 * pi * time/12)) -0.423187   0.097590  -4.336 1.45e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
## Null deviance: 343.00  on 167  degrees of freedom
```

```
## Residual deviance: 310.72  on 164  degrees of freedom
## AIC: 575.77
##
## Number of Fisher Scoring iterations: 5
```

```
plot(1970 + ((uspolio$time - 1)/12), uspolio$cases, type="h")
lines(1970 + ((uspolio$time - 1)/12), polio1.glm$fitted, col=2)
```



We can also add six-monthly components into the model.

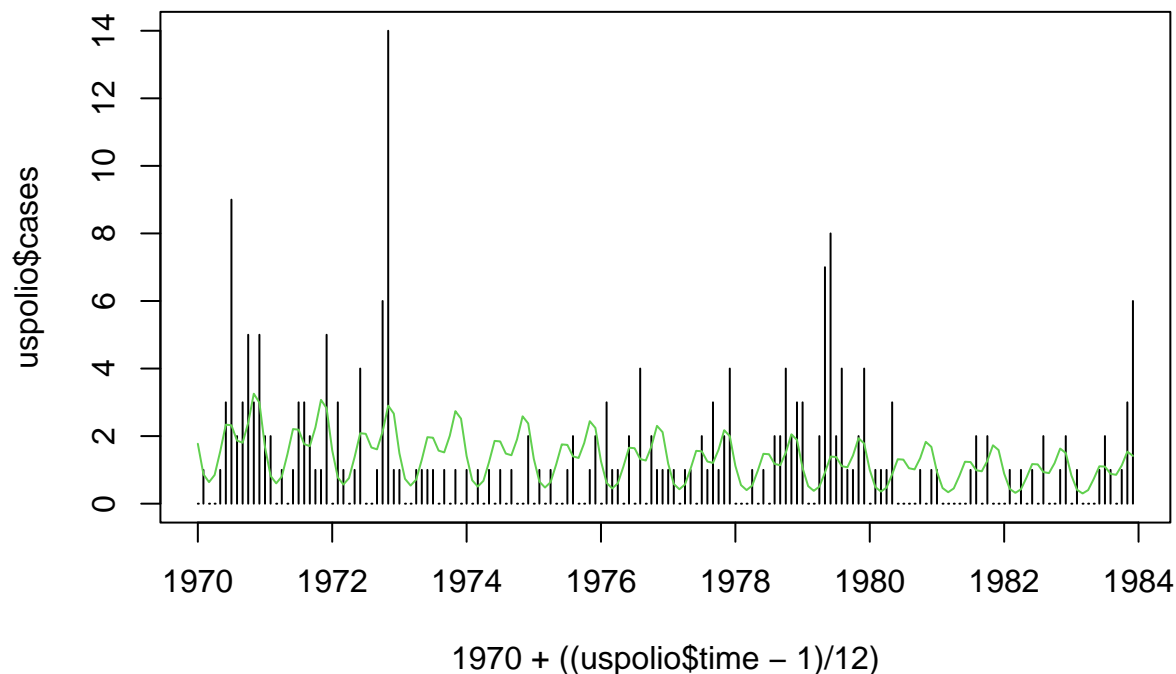
```
# Poisson model with linear trend and seasonal (annual + six-monthly) components
polio2.glm <- glm(cases ~ time + I(cos(2*pi*time/12)) + I(sin(2*pi*time/12))
+ I(cos(2*pi*time/6)) + I(sin(2*pi*time/6)), family=poisson(link=log),
data=uspolio)

summary(polio2.glm)
```

```
##
## Call:
## glm(formula = cases ~ time + I(cos(2 * pi * time/12)) + I(sin(2 *
##   pi * time/12)) + I(cos(2 * pi * time/6)) + I(sin(2 * pi *
##   time/6)), family = poisson(link = log), data = uspolio)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.557241   0.127303   4.377 1.20e-05 ***
## time          -0.004799   0.001403  -3.421 0.000625 ***
## I(cos(2 * pi * time/12))  0.137132   0.089479   1.533 0.125384
## I(sin(2 * pi * time/12)) -0.534985   0.115476  -4.633 3.61e-06 ***
## I(cos(2 * pi * time/6))   0.458797   0.101467   4.522 6.14e-06 ***
```

```
## I(sin(2 * pi * time/6)) -0.069627  0.098123 -0.710 0.477957
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 343.00  on 167  degrees of freedom
## Residual deviance: 288.85  on 162  degrees of freedom
## AIC: 557.9
##
## Number of Fisher Scoring iterations: 5
```

```
plot(1970 + ((uspolio$time - 1)/12), uspolio$cases, type="h")
lines(1970 + ((uspolio$time - 1)/12), polio2.glm$fitted, col=3)
```

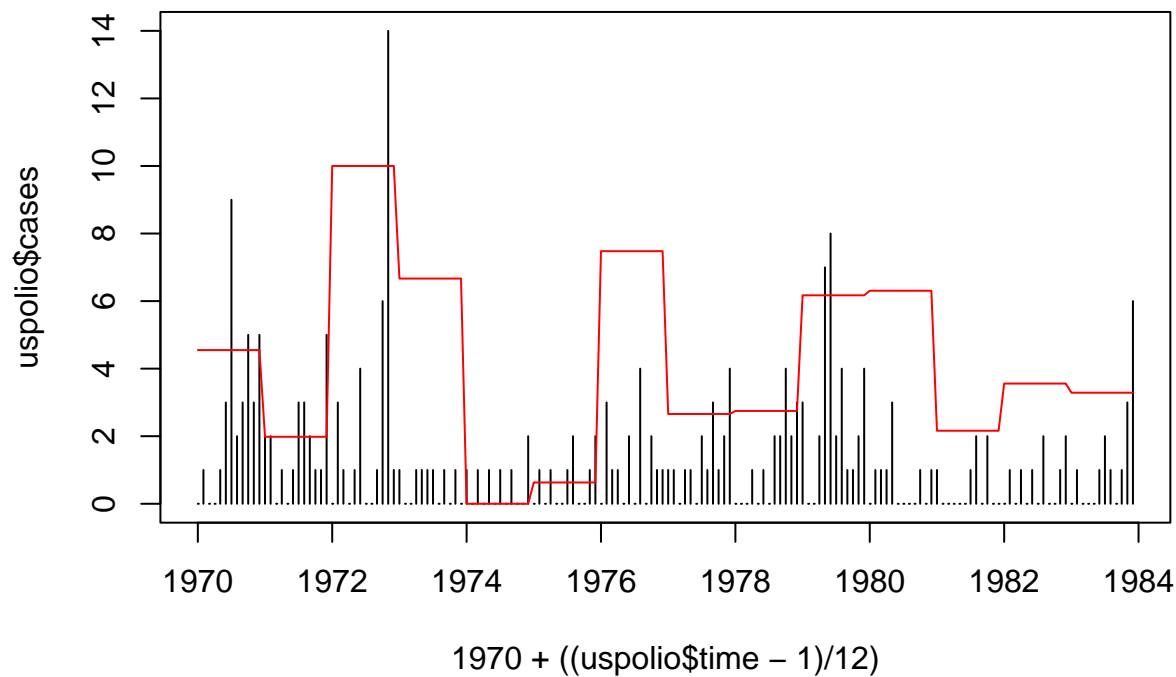


Assuming we have annual temperature data over the 14 years, we can add them into the model to investigate their effects.

```
# Average annual temperature data over the 14 years
temp_data <- rep(c(5.195, 5.138, 5.316, 5.242, 5.094, 5.108, 5.260, 5.153,
                  5.155, 5.231, 5.234, 5.142, 5.173, 5.167), each=12)

# Scale the data so that it plots nicely
scaled_temp = 10 * (temp_data - min(temp_data))/(max(temp_data) - min(temp_data))
uspolio$temp = scaled_temp

# Plot temperature data against cases data to see interest
plot(1970 + ((uspolio$time - 1)/12), uspolio$cases, type="h")
lines(1970 + ((uspolio$time - 1)/12), uspolio$temp, col="red")
```



Poisson GLM with temperature data.

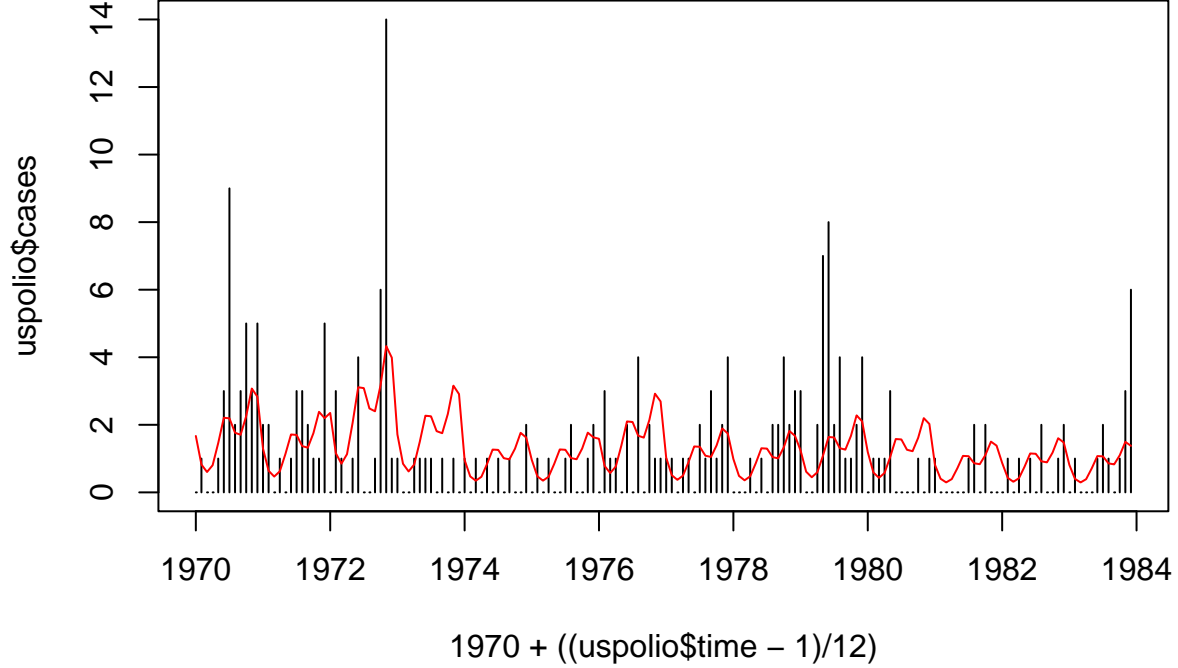
```
# Poisson model with additional temperature covariate
polio3.glm <- glm(cases ~ time + temp + I(cos(2*pi*time/12)) + I(sin(2*pi*time/12))
+ I(cos(2*pi*time/6)) + I(sin(2*pi*time/6)) , family=poisson(link=log),
data=uspolio)

summary(polio3.glm)
```

```
##
## Call:
## glm(formula = cases ~ time + temp + I(cos(2 * pi * time/12)) +
##      I(sin(2 * pi * time/12)) + I(cos(2 * pi * time/6)) + I(sin(2 *
##      pi * time/6)), family = poisson(link = log), data = uspolio)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.129643   0.186352   0.696 0.486623
## time          -0.003972   0.001439  -2.761 0.005770 **
## temp           0.080308   0.023139   3.471 0.000519 ***
## I(cos(2 * pi * time/12))  0.136094   0.089489   1.521 0.128314
## I(sin(2 * pi * time/12)) -0.531668   0.115466  -4.605 4.13e-06 ***
## I(cos(2 * pi * time/6))   0.457487   0.101435   4.510 6.48e-06 ***
## I(sin(2 * pi * time/6))  -0.068345   0.098149  -0.696 0.486218
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
```

```
##      Null deviance: 343.00  on 167  degrees of freedom
## Residual deviance: 276.84  on 161  degrees of freedom
## AIC: 547.88
##
## Number of Fisher Scoring iterations: 5
```

```
plot(1970 + ((uspolio$time - 1)/12), uspolio$cases, type="h")
lines(1970 + ((uspolio$time - 1)/12), polio3.glm$fitted, col="red")
```



2.10 Estimation of Dispersion Parameter ϕ

Because the dispersion ϕ cancels from the score equation $\mathbf{S}(\hat{\beta}) = 0$, there is no need to estimate ϕ in order to estimate β . However, $\text{Var}[\hat{\beta}]$ does depend on ϕ , as one might expect. Thus, if necessary or of interest, ϕ can be estimated via:

$$\hat{\phi} = \frac{1}{n - p} \sum_i m_i \frac{(y_i - \hat{\mu}_i)^2}{\mathcal{V}(\hat{\mu}_i)} \quad (2.88)$$

where p is the number of parameters of the model. The motivation for the above estimation is that:

$$\text{Var}[y_i] = \text{E}[(y_i - \mu_i)^2] = \phi_i \mathcal{V}(\mu_i) = \frac{\phi}{m_i} \mathcal{V}(\mu_i), \quad (2.89)$$

which can be rearranged to

$$\phi = \frac{m_i}{\mathcal{V}(\mu_i)} \text{E}[(y_i - \mu_i)^2] = \text{E} \left[m_i \frac{(y_i - \mu_i)^2}{\mathcal{V}(\mu_i)} \right]. \quad (2.90)$$

Thus, after estimating $\hat{\beta}$, we can use its value and Equation (2.88) to estimate $\hat{\phi}$.

2.10.1 Special Cases

2.10.1.1 Gaussian

When $Y|\beta, x \sim \mathcal{N}(\mu, \sigma^2)$ with $m_i = 1$, we have $\mathcal{V}(\mu_i) = 1$ and thus,

$$\hat{\phi} = \frac{1}{n-p} \sum_i (y_i - \hat{\mu}_i)^2 = \hat{\sigma}^2. \quad (2.91)$$

2.10.1.2 Gamma

Recall from Exercise 7.1 last term that we can parameterise the Gamma function in terms of its mean μ and variance σ^2 , and we found that $\mathcal{V}(\mu) = \mu^2$. Thus, when $Y|\beta, x \sim \text{Gamma}(\mu, \sigma^2)$, we have

$$\frac{1}{\hat{\nu}} = \hat{\phi} = \frac{1}{n-p} \sum_i m_i \frac{(y_i - \hat{\mu}_i)^2}{\hat{\mu}_i^2}. \quad (2.92)$$

2.10.2 Practical Example: Hospital Stay Data

In this example, we will use the Hospital Stay data introduced last term to fit a Gamma GLM and estimate its dispersion parameter.

```
library(npmlreg)
data(hosp)

# Fit the GLM and print the summary
hosp.glm <- glm(duration ~ age + temp1, data=hosp, family=Gamma(link=log))
summary(hosp.glm)

##
## Call:
## glm(formula = duration ~ age + temp1, family = Gamma(link = log),
##      data = hosp)
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -28.654096  16.621018  -1.724   0.0987 .
## age          0.014900   0.005698   2.615   0.0158 *
## temp1        0.306624   0.168141   1.824   0.0818 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Gamma family taken to be 0.2690233)
##
##      Null deviance: 8.1722  on 24  degrees of freedom
```



```
## Residual deviance: 5.7849 on 22 degrees of freedom
## AIC: 142.73
##
## Number of Fisher Scoring iterations: 6
# From the summary, note the line:
# (Dispersion parameter for Gamma family taken to be 0.2690233)

# Compute by hand
1/(hosp.glm$df.res)*sum((hosp$duration-hosp.glm$fitted)^2/(hosp.glm$fitted^2))

## [1] 0.2690233
```

2.11 Asymptotic Properties of $\hat{\beta}$

In our context, *asymptotic* means that $M = \sum_{i=1}^n m_i \rightarrow \infty$. This could be because $n \rightarrow \infty$, or because the $m_i \rightarrow \infty$, or a combination of both.

Let us denote the true value of β by β^* . In the following, we assume consistency of $\hat{\beta}$, i.e., $\hat{\beta}$ converges in probability to β^* , meaning that $P(\|\hat{\beta} - \beta^*\| \geq \varepsilon) \rightarrow 0$ as $n \rightarrow \infty$. We will denote this by $\hat{\beta} \stackrel{a}{=} \beta^*$. We will also abuse this notation to mean “tends to asymptotically” for expectations, i.e., if we write $E[Z] \stackrel{a}{=} z$, that means $E[Z] \xrightarrow{n \rightarrow \infty} z$.

From the consistency assumption, $\hat{\beta}$ will be close to β^* asymptotically, and we can expand \mathbf{S} around it:

$$\mathbf{S}(\hat{\beta}) = 0 \stackrel{a}{=} \mathbf{S}(\beta^*) + \frac{\partial \mathbf{S}(\beta^*)}{\partial \beta^T} (\hat{\beta} - \beta^*) \quad (2.93)$$

$$= \mathbf{S}(\beta^*) - \mathbf{F}_{\text{obs}}(\beta^*) (\hat{\beta} - \beta^*) \quad (2.94)$$

or equivalently,

$$\hat{\beta} - \beta^* \stackrel{a}{=} \mathbf{F}_{\text{obs}}(\beta^*)^{-1} \mathbf{S}(\beta^*). \quad (2.95)$$

2.11.1 Fisher Scoring

In Section 2.8, we stated that we often use the (expected) Fisher Information in place of the Observed Fisher Information (known as the Fisher Scoring method). Doing so in the context of asymptotic arguments is acceptable. We can roughly see this as follows. For any β ,

$$\frac{1}{n} \mathbf{F}_{\text{obs}}(\beta) = -\frac{1}{n} \frac{\partial l}{\partial \beta \partial \beta^T}(\beta) = -\frac{1}{n} \sum_{i=1}^n \frac{\partial l_i}{\partial \beta \partial \beta^T}(\beta) \rightarrow -E \left[\frac{\partial l_1}{\partial \beta \partial \beta^T}(\beta) \right] = F_1(\beta) \quad (2.96)$$

where $F_1(\beta)$ is the expected Fisher Information for a sample of size 1 and we are using the law of large numbers as $n \rightarrow \infty$ here. It can be shown (see exercise section) that $\mathbf{F}(\beta) = nF_1(\beta)$, thus justifying use of $\mathbf{F}_{\text{obs}}(\beta) \stackrel{a}{=} \mathbf{F}(\beta)$ in the forthcoming asymptotic arguments.

2.11.2 Expectation

From Equation (2.95), we have:

$$\hat{\beta} - \beta^* \stackrel{a}{=} \mathbf{F}_{\text{obs}}(\beta^*)^{-1} \mathbf{S}(\beta^*) \stackrel{a}{=} \mathbf{F}(\beta^*)^{-1} \mathbf{S}(\beta^*). \quad (2.97)$$

Because convergence in probability implies convergence in distribution, this in turn implies that

$$E[\hat{\beta} - \beta^*] \stackrel{a}{=} \mathbf{F}(\beta^*)^{-1} E[\mathbf{S}(\beta^*)] = 0. \quad (2.98)$$

In other words, $\hat{\beta}$ is asymptotically unbiased.

2.11.3 Variance

Since $E[\hat{\beta} - \beta^*] \stackrel{a}{=} 0$, we have that

$$\text{Var}[\hat{\beta} - \beta^*] \stackrel{a}{=} E[(\hat{\beta} - \beta^*)(\hat{\beta} - \beta^*)^T] \quad (2.99)$$

$$\stackrel{a}{=} E[\mathbf{F}(\beta^*)^{-1} \mathbf{S}(\beta^*) \mathbf{S}(\beta^*)^T \mathbf{F}(\beta^*)^{-T}] \quad (2.100)$$

$$= \mathbf{F}(\beta^*)^{-1} E[\mathbf{S}(\beta^*) \mathbf{S}(\beta^*)^T] \mathbf{F}(\beta^*)^{-T} \quad (2.101)$$

$$= \mathbf{F}(\beta^*)^{-1} \text{Var}[\mathbf{S}(\beta^*)] \mathbf{F}(\beta^*)^{-T} \quad (2.102)$$

$$= \mathbf{F}(\beta^*)^{-1} \quad (2.103)$$

where we have used symmetry of \mathbf{F} and the fact that $\mathbf{F}(\beta^*) = \text{Var}[\mathbf{S}(\beta^*)]$.

Thus,

$$\text{Var}[\hat{\beta}] = \text{Var}[\hat{\beta} - \beta^*] \stackrel{a}{=} \mathbf{F}(\beta^*)^{-1}. \quad (2.104)$$

2.11.4 Asymptotic Normality

The following is a sketch of the argument of asymptotic normality for $\hat{\beta} - \beta^*$, i.e., $\hat{\beta} - \beta^*$ converges asymptotically to a normal distribution. We start from

$$\mathbf{S}(\beta) = \sum_i \mathbf{S}_i(\beta) \quad (2.105)$$

where $\mathbf{S}_i(\beta)$ is defined in Section 2.6. This is a sum of independent random variables with zero mean and finite variance. As the number of terms in the sum tends to infinity, then under a certain condition, the distribution of the sum converges in distribution to a normal distribution. Since $E[\mathbf{S}(\beta)] = 0$ and $\text{Var}[\mathbf{S}(\beta)] = \mathbf{F}(\beta)$, we have:

$$\mathbf{S}(\beta) \stackrel{a}{\sim} \mathcal{N}(0, \mathbf{F}(\beta)). \quad (2.106)$$

Hence,

$$\hat{\beta} - \beta^* \stackrel{a}{=} \mathbf{F}(\beta^*)^{-1} \mathbf{S}(\beta^*) \stackrel{a}{\sim} \mathcal{N}(0, \mathbf{F}(\beta^*)^{-1} \mathbf{F}(\beta^*) \mathbf{F}(\beta^*)^{-T}). \quad (2.107)$$

Since \mathbf{F} is symmetric and convergence in probability implies convergence in distribution, we have:

$$\hat{\beta} - \beta^* \stackrel{a}{\sim} \mathcal{N}(0, \mathbf{F}(\beta^*)^{-1}). \quad (2.108)$$

This also implies that the square of Mahalanobis distance between $\hat{\beta}$ and β^* is asymptotically chi-square distributed:

$$(\hat{\beta} - \beta^*)^T \mathbf{F}(\beta^*) (\hat{\beta} - \beta^*) \stackrel{a}{\sim} \chi^2(p) \quad (2.109)$$

where p is the number of parameters.

2.11.5 Closing The Circle

At the beginning of this section, we assumed that $\hat{\beta}$ converges in probability to β^* . We may want to justify that this assumption is reasonable. Note that under some regularity conditions,

$$\mathbf{F}(\beta)^{-1} = \left(\sum_i m_i \dots \right)^{-1} \rightarrow 0 \quad (2.110)$$

as $M \rightarrow \infty$. Thus $\hat{\beta}$ converges in distribution to a constant random variable, which means that it converges in probability too. This is what we were assuming.

Equations (2.104), (2.108), and (2.109) remain valid when $\mathbf{F}(\beta^*)$ is replaced by $\mathbf{F}(\hat{\beta})$.

2.11.6 Next Step

Now that we have seen how to estimate the parameters and some of their sampling properties (asymptotically), we can move on to use these estimates to make inferences: predictions about new values and confidence intervals.

2.12 Exercises

1. Show that $\mathbf{F}(\beta) = nF_1(\beta)$, where $\mathbf{F}(\beta)$ and $F_1(\beta)$ are the Fisher Information on datasets of size n and 1 respectively.
2. Using the notation and argument in Section 2.11, show that $\mathbf{F}(\beta^*) \stackrel{a}{=} \mathbf{F}(\hat{\beta})$ and thus we can replace $\mathbf{F}(\beta^*)$ by $\mathbf{F}(\hat{\beta})$ in Equations (2.104), (2.108), and (2.109).

Bibliography

P. J. Green. Iteratively reweighted least squares for maximum likelihood estimation, and some robust and resistant alternatives. *JRSSB*, 46(2):149–192, 1984.

John A Nelder and Robert WM Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 135(3):370–384, 1972.