

NMST432 Advanced Regression Models

Course notes

Michal Kulich

Last modified on May 22, 2020.



matfyz

Department of Probability and Mathematical Statistics
Faculty of Mathematics and Physics, Charles University

These course notes contain an overview of notation, definitions, theorems and comments covered by the course “NMST432 Advanced regression models”, which is a part of the curriculum of the Master’s program “Probability, Mathematical Statistics and Econometrics”.

This document undergoes continuing development. The author will appreciate notifications by the reader of potential typos or misprints.

Michal Kulich
kulich@karlin.mff.cuni.cz

In Karlín on May 22, 2020

Contents

1	Review of Linear Regression	6
1.1	Definition and Assumptions	6
1.2	Estimation	7
1.3	Normal Linear Regression	8
1.4	Asymptotic Properties of the LSE	9
1.5	Transformed Response	10
2	Generalized Linear Model – Theory	11
2.1	Exponential Family	11
2.1.1	Parametrization, moments	11
2.1.2	Maximum likelihood estimator of the canonical parameter	14
2.2	Definition of the Generalized Linear Model	15
2.3	Maximum Likelihood Estimation in the GLM	19
2.4	Algorithm for Fitting the GLM	22
2.5	Estimation of the Dispersion Parameter	23
2.6	Deviance	24
2.7	Asymptotic Results	25
2.8	Diagnostic Methods for the GLM	29
2.8.1	Pearson residuals	29
2.8.2	Leverages	30
2.8.3	Standardized Pearson residuals	30
2.8.4	Deviance residuals	30
2.8.5	Standardized deviance residuals	30
2.8.6	Cook’s distance	31
2.8.7	Residual plots	31
2.8.8	Diagnostics of the link function	31
2.9	Model-building strategies	31
3	Generalized Linear Model for Discrete Responses	34
3.1	Analysis of Binary Data	34
3.1.1	Alternative vs. binomial data	34
3.1.2	Link functions for binary data	35
3.1.3	Binary data likelihood	36
3.1.4	Threshold analysis by probit regression	37
3.1.5	Logistic regression	38

3.2	Analysis of Poisson Count Data	42
3.2.1	Poisson loglinear model	42
3.2.2	Modelling Poisson process intensity	44
3.3	Loglinear Models for Contingency Tables	45
3.3.1	Two-way contingency table	45
3.3.2	Distributions of observed counts	45
3.3.3	Loglinear models for two-way tables	48
3.3.4	Testing independence in a two-way table	50
3.3.5	Loglinear models for three-way tables	50
3.3.6	Loglinear models for multi-way tables	57
3.3.7	Equivalence of loglinear and logistic models	57
4	Extensions of Generalized Linear Models	60
4.1	Quasi-likelihood and Overdispersion	60
4.1.1	Overdispersion in binomial data	60
4.1.2	Overdispersion in Poisson data	60
4.1.3	Quasi-likelihood	61
4.2	Sandwich Variance Estimation in the GLM	63
4.2.1	Behavior of the MLE under a misspecified model	63
4.2.2	Applications to the GLM	65
5	Generalized Estimating Equations	67
5.1	Group-Dependent Data	67
5.2	Estimation of Regression Parameters by Generalized Estimating Equations	68
5.3	Correlation Structures	69
5.3.1	Working independence	70
5.3.2	Parametrized correlations	70
5.3.3	Joint estimation of mean and correlation structures	72
5.3.4	Summary of GEE methods	72
6	Linear Mixed Effects Models	73
6.1	Introduction	73
6.1.1	One-way ANOVA	73
6.1.2	One-way ANOVA with random effects	75
6.1.3	Two-way ANOVA with random effects	77
6.1.4	Random intercept and slope	79
6.2	Definition of Linear Mixed Effects Model	80
6.2.1	Single-level LME model	80
6.2.2	Multi-level LME model	82
6.3	Parameter Estimation	82
6.3.1	Marginal likelihood	82
6.3.2	Henderson's mixed model equations	83
6.3.3	Maximum likelihood estimation of variance parameters	85

6.3.4	Restricted maximum likelihood (REML) estimation of variance parameters	87
6.3.5	Comparison of REML versus ML estimators	90
6.4	Hypothesis Testing and Confidence Intervals	91
6.4.1	Asymptotic approach based on MLE theory	91
6.4.2	Likelihood ratio tests	92
6.4.3	t tests and F tests for fixed effects	94
6.5	Extended Linear Mixed Effects Model	96
6.5.1	Introduction	96
6.5.2	Parameter estimation	96
6.5.3	Generalized least squares	97
6.5.4	Decomposing variance structure	98
6.6	Comparison of LME and GEE Approaches	99
7	Generalized Linear Mixed Models	100
7.1	Model and Assumptions	100
7.2	Parameter Estimation	101
7.3	Interpretation	103
7.4	Comparison of GLMM vs. GEE models	104
	Bibliography	105
	Index	107

1 Review of Linear Regression

1.1 Definition and Assumptions

Consider n independent copies of random vectors (Y_i, \mathbf{X}_i) , $i = 1, \dots, n$. Each \mathbf{X}_i has $p < n$ components (X_{i1}, \dots, X_{ip}) .

Note.

- Y_i is called *the response*^{*}. The components of \mathbf{X}_i are called *covariates* (explanatory variables, predictors, regressors)[†].
- The covariate X_{i1} is usually taken as 1.
- In certain applications, the covariates can be fixed quantities rather than random variables. Throughout this course, we will consider covariates random. Extensions to fixed covariates usually hold with some additional conditions but the proofs require more effort.

Notation. Let $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ and

$$\mathbb{X} = \begin{pmatrix} \mathbf{X}_1^\top \\ \mathbf{X}_2^\top \\ \vdots \\ \mathbf{X}_n^\top \end{pmatrix}.$$

The n by p matrix \mathbb{X} is called *the regression matrix*[‡]. We assume $r(\mathbb{X}) = p$ (full rank).

Definition 1.1. The data (Y_i, \mathbf{X}_i) satisfy the linear regression model if

$$Y_i = \mathbf{X}_i^\top \boldsymbol{\beta}_0 + \varepsilon_i,$$

where $\varepsilon_1, \dots, \varepsilon_n$ are independent, $\mathbf{E}[\varepsilon_i | \mathbf{X}_i] = 0$, and $\text{var}[\varepsilon_i | \mathbf{X}_i] = \sigma^2$.

Note. $\boldsymbol{\beta}_0 = (\beta_{01}, \dots, \beta_{0p})^\top$ is a vector of unknown parameters (*regression coefficients*[§]), random variables ε_i are called *error terms (disturbances)*[¶], σ^2 is called *residual variance*^{||}.

Note. Another convenient formulation of the model is based on conditional moments (avoids the introduction of error terms). The linear regression model holds if and only if

^{*} Český odezva [†] Český regresory, nezávisle proměnné, vysvětlující veličiny, prediktory, kovariáty
[‡] Český regresní matice [§] Český regresní koeficienty [¶] Český chybové členy ^{||} Český residuální rozptyl

- Y_1, \dots, Y_n are independent
- $\mathbb{E}[Y_i | \mathbf{X}_i] = \mathbf{X}_i^\top \boldsymbol{\beta}_0$
- $\text{var}[Y_i | \mathbf{X}_i] = \sigma^2$

The model specifies the first two conditional moments of Y_i given \mathbf{X}_i .

Note. We will use the notation \mathbb{E} , var for conditional moments given \mathbf{X}_i . The symbol \mathbb{E}_X will be used for unconditional expectation over the distribution of \mathbf{X}_i .

Note. The regression parameters express the influence of \mathbf{X}_i on $\mathbb{E}Y_i$. Assuming that $X_{i1} = 1$,

$$\beta_{01} = \mathbb{E}[Y_i | X_{i2} = 0, X_{i3} = 0, \dots, X_{ip} = 0]$$

and

$$\beta_{02} = \mathbb{E}[Y_i | X_{i2} = x + 1, X_{i3} = x_3, \dots, X_{ip} = x_p] - \mathbb{E}[Y_i | X_{i2} = x, X_{i3} = x_3, \dots, X_{ip} = x_p]$$

1.2 Estimation

The regression coefficients $\boldsymbol{\beta}_0$ are estimated by *the least squares estimator* (LSE) $\hat{\boldsymbol{\beta}}$ that minimizes the sum of squares

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^n (Y_i - \mathbf{X}_i^\top \boldsymbol{\beta})^2 = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} (\mathbf{Y} - \mathbb{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbb{X}\boldsymbol{\beta}),$$

i.e., solves the system of *normal equations*

$$\sum_{i=1}^n \mathbf{X}_i (Y_i - \mathbf{X}_i^\top \hat{\boldsymbol{\beta}}) = \mathbf{0}.$$

Because \mathbb{X} is of full rank, the single solution to the system is

$$\hat{\boldsymbol{\beta}} = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{Y}.$$

Note.

- $\mathbb{E} \hat{\boldsymbol{\beta}} = \boldsymbol{\beta}_0$ (unbiased), $\text{var} \hat{\boldsymbol{\beta}} = \sigma^2 (\mathbb{X}^\top \mathbb{X})^{-1}$.
- The vector

$$\hat{\mathbf{Y}} = \mathbb{X} \hat{\boldsymbol{\beta}} = \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{Y} = \mathbb{H} \mathbf{Y}$$

is called *the vector of fitted values**.

- The projection matrix $\mathbb{H} \stackrel{\text{df}}{=} \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top$ is idempotent, with rank p . It satisfies $\mathbb{H} \mathbb{X} = \mathbb{X}$. The matrix $\mathbb{I}_n - \mathbb{H}$ is also idempotent with rank $n - p$, and satisfies $(\mathbb{I}_n - \mathbb{H}) \mathbb{X} = \mathbf{0}$.

* Česky vektor odhadnutých (vyrovnaných) hodnot

- $E\hat{Y} = \mathbb{X}\beta_0$, $\text{var}\hat{Y} = \sigma^2\mathbb{H}$.
- The random vector $\mathbf{u} \stackrel{\text{df}}{=} \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbb{I}_n - \mathbb{H})\mathbf{Y}$ is called *the vector of residuals**. It satisfies $E\mathbf{u} = \mathbf{0}$, $\text{var}\mathbf{u} = \sigma^2(\mathbb{I}_n - \mathbb{H})$.
- The random variable

$$SS_e \stackrel{\text{df}}{=} \mathbf{u}^\top \mathbf{u} = \sum_{i=1}^n (Y_i - \mathbf{X}_i^\top \hat{\boldsymbol{\beta}})^2 = \mathbf{Y}^\top (\mathbb{I}_n - \mathbb{H}) \mathbf{Y}$$

is called the *residual sum of squares*†. Because $E SS_e = (n - p)\sigma^2$, we obtain an unbiased estimator of residual variance as $\hat{\sigma}^2 = SS_e/(n - p)$.

1.3 Normal Linear Regression

For normally distributed errors, additional useful properties can be derived. Assume now that $\boldsymbol{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2\mathbb{I}_n)$.

Proposition 1.1. *Under normality,*

- (i) $\mathbf{Y} \sim N_n(\mathbb{X}\boldsymbol{\beta}_0, \sigma^2\mathbb{I}_n)$
- (ii) $\hat{\boldsymbol{\beta}} \sim N_p(\boldsymbol{\beta}_0, \sigma^2(\mathbb{X}^\top\mathbb{X})^{-1})$
- (iii) $\hat{\mathbf{Y}} \sim N_n(\mathbb{X}\boldsymbol{\beta}_0, \sigma^2\mathbb{H})$
- (iv) $\mathbf{u} \sim N_n(\mathbf{0}, \sigma^2(\mathbb{I}_n - \mathbb{H}))$
- (v) $SS_e/\sigma^2 \sim \chi_{n-p}^2$
- (vi) $\hat{\boldsymbol{\beta}}$ and SS_e are independent
- (vii) Let \mathbf{c} be any non-zero p -vector of real constants. Then

$$\frac{\mathbf{c}^\top \hat{\boldsymbol{\beta}} - \mathbf{c}^\top \boldsymbol{\beta}_0}{\sqrt{\hat{\sigma}^2 \mathbf{c}^\top (\mathbb{X}^\top\mathbb{X})^{-1} \mathbf{c}}} \sim t_{n-p}$$

- (viii) Assume the model $\mathbf{Y} = \mathbb{X}\boldsymbol{\beta}_0 + \boldsymbol{\varepsilon}$, where $\mathbb{X} = (\mathbb{X}_A | \mathbb{X}_B)$ and $\boldsymbol{\beta}_0 = (\boldsymbol{\beta}_A^\top, \boldsymbol{\beta}_B^\top)^\top$, $\boldsymbol{\beta}_B \in \mathbb{R}_m$, $\boldsymbol{\beta}_A \in \mathbb{R}^{p-m}$, and introduce the submodel $\mathbf{Y} = \mathbb{X}_A\boldsymbol{\beta}_A + \boldsymbol{\varepsilon}'$. Let SS_e and SS_h be the residual sums of squares in the model and submodel, respectively. If the submodel is true ($H_0 : \boldsymbol{\beta}_B = \mathbf{0}$ holds) then

$$F = \frac{n-p}{m} \frac{SS_h - SS_e}{SS_e} \sim F_{m, n-p}. \quad (1.1)$$

It can be also shown that, under normality, $\hat{\boldsymbol{\beta}}$ is the best linear unbiased estimator and the maximum likelihood estimator, so it possesses optimality properties.

* Český vektor residuí † Český residuální součet čtverců

1.4 Asymptotic Properties of the LSE

Let (Y_i, \mathbf{X}_i) , $i = 1, \dots, n$, be iid. Assume Definition 1.1 (without normality). Denote $\mathbb{D}_X = \mathbb{E}_X \mathbf{X}_i \mathbf{X}_i^\top$.

Proposition 1.2. *Let \mathbb{D}_X be a finite regular matrix. Then*

- (i) $\hat{\boldsymbol{\beta}} \xrightarrow{P} \boldsymbol{\beta}_0$ as $n \rightarrow \infty$,
- (ii) $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \xrightarrow{D} \mathbf{N}_p(\mathbf{0}, \sigma^2 \mathbb{D}_X^{-1})$ as $n \rightarrow \infty$.

Proposition 1.2(ii) is an asymptotic restatement of Proposition 1.1(ii). Other parts of Proposition 1.1 also hold asymptotically even if the data are not normal.

Now relax the assumption of equal variance: assume only $\mathbb{E}[Y_i | \mathbf{X}_i] = \mathbf{X}_i^\top \boldsymbol{\beta}_0$. Let $\text{var}[Y_i | \mathbf{X}_i] = \sigma^2(\mathbf{X}_i)$ be stochastically bounded (finite expectation follows). Denote $\mathbb{V}_X = \mathbb{E}_X \sigma^2(\mathbf{X}_i) \mathbf{X}_i \mathbf{X}_i^\top$.

Proposition 1.3. *Let \mathbb{V}_X be finite and \mathbb{D}_X be finite and regular. Then*

- (i) $\hat{\boldsymbol{\beta}} \xrightarrow{P} \boldsymbol{\beta}_0$ as $n \rightarrow \infty$,
- (ii) $\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \xrightarrow{D} \mathbf{N}_p(\mathbf{0}, \mathbb{D}_X^{-1} \mathbb{V}_X \mathbb{D}_X^{-1})$ as $n \rightarrow \infty$.

When equal variances hold, $\mathbb{V}_X = \sigma^2 \mathbb{D}_X$ and the result in Proposition 1.3(ii) transforms into the result in Proposition 1.2(ii).

Consistent estimates of \mathbb{D}_X and \mathbb{V}_X are

$$\hat{\mathbb{D}}_n = \frac{1}{n} \mathbb{X}^\top \mathbb{X}$$

and

$$\hat{\mathbb{V}}_n = \frac{1}{n} \mathbb{X}^\top \text{diag}(u_i^2) \mathbb{X}.$$

So, if both normality and homoskedasticity are in doubt, one can use the OLS estimator $\hat{\boldsymbol{\beta}}$ with variance

$$(\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \text{diag}(u_i^2) \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1}$$

in place of the usual $\hat{\sigma}^2 (\mathbb{X}^\top \mathbb{X})^{-1}$. This is called *the sandwich estimator**, or, in the econometric context, *White estimator†* (White 1980).

Many variants and improvements of this estimator have been proposed in the literature.

* Český sendvičový odhad † Český Whiteův odhad

1.5 Transformed Response

Recall the linear model $E[Y_i | \mathbf{X}_i] = \mathbf{X}_i^\top \boldsymbol{\beta}_0$ with $\text{var}[Y_i | \mathbf{X}_i] = \sigma^2$. The regression parameters can be interpreted as

$$\beta_{01} = E[Y_i | X_{i2} = 0, X_{i3} = 0, \dots, X_{ip} = 0]$$

and

$$\beta_{02} = E[Y_i | X_{i2} = x + 1, X_{i3} = x_3, \dots, X_{ip} = x_p] - E[Y_i | X_{i2} = x, X_{i3} = x_3, \dots, X_{ip} = x_p].$$

When the response is non-normal, the common practice is to specify a linear model on a transformed response. Let g be some monotone function. The transformed model is

$$g(Y_i) = \mathbf{X}_i^\top \boldsymbol{\beta}_0 + \varepsilon_i$$

or $E[g(Y_i) | \mathbf{X}_i] = \mathbf{X}_i^\top \boldsymbol{\beta}_0$ with $\text{var}[g(Y_i) | \mathbf{X}_i] = \sigma^2$. The induced model for Y_i is

$$Y_i = g^{-1}(\mathbf{X}_i^\top \boldsymbol{\beta}_0 + \varepsilon_i).$$

In general, the effect of the covariates on EY_i in this model cannot be expressed.

The only special case (apart from linear g) when the transformed model says anything useful about $E[Y_i | \mathbf{X}_i]$ is the log transform. From

$$\log Y_i = \mathbf{X}_i^\top \boldsymbol{\beta}_0 + \varepsilon_i$$

we get a multiplicative model

$$Y_i = e^{\mathbf{X}_i^\top \boldsymbol{\beta}_0} \varepsilon_i^*,$$

where $\varepsilon_i^* = e^{\varepsilon_i}$, $E \varepsilon_i^* = \mu_\varepsilon > 1$, $\text{var} \varepsilon_i^* = \sigma_\varepsilon^2$. Then

$$\begin{aligned} E[Y_i | \mathbf{X}_i] &= \exp\{\log \mu_\varepsilon + \mathbf{X}_i^\top \boldsymbol{\beta}_0\}, \\ \text{var}[Y_i | \mathbf{X}_i] &= \sigma_\varepsilon^2 (\exp\{\mathbf{X}_i^\top \boldsymbol{\beta}_0\})^2. \end{aligned}$$

While β_{01} (the intercept) does not have useful interpretation, the other parameters express multiplicative effects of X_{i2}, \dots, X_{ip} on EY_i . E.g.,

$$e^{\beta_{02}} = \frac{E[Y_i | X_{i2} = x + 1, X_{i3} = x_3, \dots, X_{ip} = x_p]}{E[Y_i | X_{i2} = x, X_{i3} = x_3, \dots, X_{ip} = x_p]}.$$

So, $e^{\beta_{0k}}$ is the proportional increase (relative change) in EY_i after a unit change in X_{ik} .

The problem with the interpretation of the transformed linear model is serious when the primary task is to estimate the effect of \mathbf{X}_i on EY_i . If the goal is to predict Y_i from \mathbf{X}_i , transformations can still be useful even if the interpretation of the parameters is lost.

2 Generalized Linear Model – Theory

The generalized linear model extends the normal linear model in two ways:

- allows a wider choice of distributions for Y_i (distributions from exponential family);
- allows a more general relationship between $E Y_i$ and $\mathbf{X}_i^T \boldsymbol{\beta}_0$.

2.1 Exponential Family

2.1.1 Parametrization, moments

Definition 2.1. If the density of a real random variable (w.r.t. some σ -finite measure μ) can be written in the form

$$f(x; \theta, \varphi) = \exp\left\{\frac{x\theta - b(\theta)}{\varphi} + c(x, \varphi)\right\}, \quad (2.1)$$

where

- θ is called *the canonical parameter*^{*};
- $\varphi \in (0, \infty)$ is called *the dispersion parameter*[†];
- b and c are some real functions;

then the distribution belongs to *the exponential family* of distributions[‡]. The expression (2.1) is called *the canonical form of the density*[§].

Examples

Normal distribution $Y \sim N(\mu, \sigma^2)$

$$\begin{aligned} f(x; \mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}, \quad \mu \in \mathbb{R}, \sigma^2 > 0, x \in \mathbb{R} \\ &= \exp\left\{\frac{x\mu - \mu^2/2}{\sigma^2} - \frac{x^2}{2\sigma^2} - \frac{1}{2}\log(2\pi\sigma^2)\right\} \\ \theta &= \mu, \quad \varphi = \sigma^2, \quad b(\theta) = \frac{\theta^2}{2}, \quad c(x, \varphi) = -\frac{x^2}{2\varphi} - \frac{1}{2}\log(2\pi\varphi) \end{aligned}$$

^{*} Český kanonický parametr [†] Český disperzní parametr [‡] Český rozdělení exponenciálního typu
[§] Český kanonický tvar hustoty

Gamma distribution $Y \sim \Gamma(a, p)$

$$\begin{aligned} f(x; a, p) &= \frac{a^p}{\Gamma(p)} x^{p-1} \exp\{-ax\}, \quad a, p > 0, \quad x > 0 \\ &= \exp\left\{\frac{-(a/p)x + \log(a/p)}{1/p} + (p-1) \log x + p \log p - \log \Gamma(p)\right\} \\ \theta &= -\frac{a}{p}, \quad \varphi = 1/p, \quad b(\theta) = -\log(-\theta) \\ c(x, \varphi) &= (1/\varphi - 1) \log x - (\log \varphi)/\varphi - \log \Gamma(1/\varphi) \end{aligned}$$

Inverse Gaussian distribution $Y \sim \text{IG}(\mu, \lambda)$

$$\begin{aligned} f(x; \mu, \lambda) &= \sqrt{\frac{\lambda}{2\pi x^3}} \exp\left\{-\frac{\lambda(x-\mu)^2}{2\mu^2 x}\right\}, \quad \mu, \lambda > 0, \quad x > 0 \\ &= \exp\left\{\frac{-x/(2\mu^2) + 1/\mu}{1/\lambda} + \frac{1}{2} \log \frac{\lambda}{2\pi x^3} - \frac{\lambda}{2x}\right\} \\ \theta &= -\frac{1}{2\mu^2}, \quad \varphi = 1/\lambda, \quad b(\theta) = -\sqrt{-2\theta}, \quad c(x, \varphi) = -\frac{1}{2} \log(2\pi x^3 \varphi) - (2x\varphi)^{-1} \end{aligned}$$

Poisson distribution $Y \sim \text{Po}(\lambda)$

$$\begin{aligned} f(x; \lambda) &= \frac{\lambda^x}{x!} \exp\{-\lambda\}, \quad \lambda > 0, \quad x = 0, 1, 2, \dots \\ &= \exp\{x \log \lambda - \lambda + \log x!\} \\ \theta &= \log \lambda, \quad \varphi = 1, \quad b(\theta) = \exp(\theta), \quad c(x, \varphi) = \log x! \end{aligned}$$

Alternative distribution $Y \sim \text{Alt}(p)$

$$\begin{aligned} f(x; p) &= p^x (1-p)^{1-x}, \quad p \in (0, 1), \quad x = 0, 1 \\ &= \exp\left\{x \log \frac{p}{1-p} + \log(1-p)\right\} \\ \theta &= \log \frac{p}{1-p}, \quad \varphi = 1, \quad b(\theta) = \log(1 + \exp\{\theta\}), \quad c(x, \varphi) = 0 \end{aligned}$$

Lemma 2.1. Let the random variable Y follow a distribution from the exponential family. Then the moment generation function $m_Y(t) \equiv \mathbb{E} e^{tY}$ of Y exists, is finite, and is equal to

$$m_Y(t) = \exp\left\{\frac{b(t\varphi + \theta) - b(\theta)}{\varphi}\right\}.$$

If $b(\theta)$ is twice continuously differentiable, $m_Y(t)$ is twice differentiable at $t = 0$.

Corollary. If $b(\theta)$ is twice continuously differentiable then Y has finite first two moments, $EY = b'(\theta)$ and $\text{var } Y = \varphi b''(\theta)$.

We will always assume that $b(\theta)$ is twice continuously differentiable so that $\text{var } Y$ is finite. Let $\mu \stackrel{\text{df}}{=} EY$.

Note. Since $\text{var } Y = \varphi b''(\theta) > 0$, b is a strictly convex function and b' is strictly increasing. Hence b' has a well-defined inverse. There exists a function $V(\mu)$ of the mean μ such that $\text{var } Y = \varphi V(\mu)$. It satisfies the equation $b''(\theta) = V(b'(\theta))$.

Definition 2.2. The function $V(\mu)$ is called *the variance function**.

Note.

- Each distribution belonging to the exponential family has a different variance function.
- Within the exponential family, the variance function determines the distribution of Y . Not every function V is a variance function of some distribution from the exponential family.

Examples

Normal distribution $Y \sim N(\mu, \sigma^2)$

$$\theta = \mu, \quad \varphi = \sigma^2, \quad b(\theta) = \frac{\theta^2}{2}$$

$$EY = b'(\theta) = \mu, \quad \text{var } Y = \varphi b''(\theta) = \varphi = \sigma^2, \quad V(\mu) = 1$$

This is the only distribution in exponential family with constant variance function, i.e., the variance is unrelated to the mean.

Gamma distribution $Y \sim \Gamma(a, p)$

$$\theta = -\frac{a}{p}, \quad \varphi = 1/p, \quad b(\theta) = -\log(-\theta)$$

$$\mu = EY = b'(\theta) = -1/\theta = p/a, \quad \text{var } Y = \varphi b''(\theta) = \varphi/\theta^2 = p/a^2, \quad V(\mu) = \mu^2$$

Inverse Gaussian distribution $Y \sim \text{IG}(\mu, \lambda)$

$$\theta = -\frac{1}{2\mu^2}, \quad \varphi = 1/\lambda, \quad b(\theta) = -\sqrt{-2\theta}$$

$$EY = b'(\theta) = 1/\sqrt{-2\theta} = \mu, \quad \text{var } Y = \varphi b''(\theta) = \varphi(-2\theta)^{-3/2} = \mu^3/\lambda, \quad V(\mu) = \mu^3$$

* Český rozptylová funkce

Poisson distribution $Y \sim \text{Po}(\lambda)$

$$\theta = \log \lambda, \quad \varphi = 1, \quad b(\theta) = \exp(\theta)$$

$$\mu = \mathbf{E}Y = b'(\theta) = \exp(\theta) = \lambda, \quad \text{var } Y = \varphi b''(\theta) = \exp(\theta) = \lambda, \quad V(\mu) = \mu$$

Alternative distribution $Y \sim \text{Alt}(p)$

$$\theta = \log \frac{p}{1-p}, \quad \varphi = 1, \quad b(\theta) = \log(1 + \exp\{\theta\}) = \log(1-p)$$

$$\mu = \mathbf{E}Y = b'(\theta) = \frac{e^\theta}{1+e^\theta} = p, \quad \text{var } Y = \varphi b''(\theta) = \frac{e^\theta}{(1+e^\theta)^2} = p(1-p),$$

$$V(\mu) = \mu(1-\mu)$$

2.1.2 Maximum likelihood estimator of the canonical parameter

Let Y_1, \dots, Y_n be a random sample from the density $f(x; \theta_0, \varphi_0)$ belonging to the exponential family, θ_0 is the true canonical parameter, φ_0 is the true dispersion parameter. Let $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$.

The likelihood is

$$L(\theta, \varphi) = \prod_{i=1}^n \exp\left\{ \frac{Y_i \theta - b(\theta)}{\varphi} + c(Y_i, \varphi) \right\},$$

The log-likelihood is

$$\ell(\theta, \varphi) = \log L(\theta, \varphi) = \sum_{i=1}^n \left[\frac{Y_i \theta - b(\theta)}{\varphi} + c(Y_i, \varphi) \right].$$

Suppose that the true dispersion parameter φ_0 is known. Then the score function for θ is

$$U(\theta | Y_i) = \frac{\partial}{\partial \theta} \log f(x; \theta, \varphi_0) = \frac{1}{\varphi_0} [Y_i - b'(\theta)].$$

Obviously, $\mathbf{E}U(\theta_0 | Y_i) = 0$. The score statistic is

$$U_n(\theta | \mathbf{Y}) = \sum_{i=1}^n U(\theta | Y_i) = \frac{1}{\varphi_0} \sum_{i=1}^n [Y_i - b'(\theta)].$$

The maximum likelihood estimator [MLE] $\hat{\theta}_n$ solves the equation $U_n(\hat{\theta}_n | \mathbf{Y}) = 0$, that is $\sum_{i=1}^n Y_i = nb'(\hat{\theta}_n)$. The solution is $\hat{\theta}_n = (b')^{-1}(\bar{Y}_n)$, where $\bar{Y}_n = n^{-1} \sum_{i=1}^n Y_i$.

The MLE is unique because b is convex, and it does not depend on the dispersion parameter φ_0 . It can be calculated even if φ_0 is unknown.

The observed information is

$$I_n(\theta | \mathbf{Y}) = -\frac{1}{n} \sum_{i=1}^n \frac{\partial U(\theta | Y_i)}{\partial \theta} = \frac{1}{\varphi_0} b''(\theta) > 0,$$

so the likelihood is strictly concave. The expected (Fisher) information is the same as the observed information,

$$I(\theta) = -\mathbb{E} \frac{\partial U(\theta | Y_i)}{\partial \theta} = \frac{1}{\varphi_0} b''(\theta).$$

It is easy to check that

$$\text{var} U(\theta_0 | Y_i) = I(\theta_0).$$

It follows from the theory of MLE that

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{D} \mathbf{N}(0, \varphi_0 [b''(\theta_0)]^{-1}). \quad (2.2)$$

Now consider the true dispersion parameter φ_0 unknown. The MLE of θ_0 is still the same, $\hat{\theta}_n = (b')^{-1}(\bar{Y}_n)$. However, what is the asymptotic distribution of $\hat{\theta}_n$ when φ_0 is unknown? In general, the asymptotic variance may change.

Calculate the joint information matrix for (θ, φ) :

$$I(\theta_0, \varphi_0) = -\mathbb{E} \frac{\partial^2 \log f(x; \theta_0, \varphi_0)}{\partial(\theta, \varphi) \partial(\theta, \varphi)^\top} = \begin{pmatrix} I_{\theta\theta} & I_{\theta\varphi} \\ I_{\theta\varphi} & I_{\varphi\varphi} \end{pmatrix} = \begin{pmatrix} b''(\theta_0)/\varphi_0 & 0 \\ 0 & I_{\varphi\varphi} \end{pmatrix}.$$

Thus, the information matrix is diagonal. It follows that the asymptotic distribution of $\hat{\theta}_n$ is given by (2.2) even if φ_0 is unknown.

We do not need φ_0 to estimate θ_0 but we need an estimate of φ_0 to estimate the asymptotic variance of θ_0 . Of course, we could use the MLE of φ_0 but it often cannot be calculated explicitly. Instead, we can use the moment estimator

$$\hat{\varphi} = \frac{S_n^2}{b''(\hat{\theta}_n)} = \frac{S_n^2}{V(\bar{Y}_n)},$$

where S_n^2 is the sample variance. Since $S_n^2 \xrightarrow{P} \text{var} Y_i = \varphi_0 b''(\theta_0)$, $\hat{\theta}_n$ is consistent and b'' is continuous, $\hat{\varphi}$ is consistent (though less efficient than the MLE).

2.2 Definition of the Generalized Linear Model

Consider n independent copies of random vectors (Y_i, \mathbf{X}_i) , $i = 1, \dots, n$, where $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^\top$. We want to express the dependence of $\mu_i \stackrel{\text{df}}{=} \mathbb{E}[Y_i | \mathbf{X}_i]$ on \mathbf{X}_i by a model that is more general than the linear model.

Definition 2.3. (Nelder and Wedderburn 1972) The data (Y_i, \mathbf{X}_i) satisfy the *generalized linear model** [GLM] if

1. Y_1, \dots, Y_n are independent and the distribution of Y_i depends on \mathbf{X}_i through regression parameters $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$.
2. the conditional density of Y_i given \mathbf{X}_i has the form

$$f(y; \theta_i, \varphi) = \exp \left\{ \frac{y\theta_i - b(\theta_i)}{\varphi} + c(y, \varphi) \right\},$$

(is of exponential type), where $b(\cdot)$ is a known twice continuously differentiable function, θ_i depends on \mathbf{X}_i and $\boldsymbol{\beta}$, $\varphi > 0$ is a known or an unknown constant.

3. θ_i depends on \mathbf{X}_i and $\boldsymbol{\beta}$ through the *linear predictor*[†] $\eta_i \stackrel{\text{df}}{=} \mathbf{X}_i^\top \boldsymbol{\beta}$.
4. There exists a known strictly monotone, twice continuously differentiable *link function*[‡] g such that $g(\mu_i) = \eta_i$.

Notation. Let $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ and define the regression matrix

$$\mathbb{X}_{n \times p} = \begin{pmatrix} \mathbf{X}_1^\top \\ \mathbf{X}_2^\top \\ \vdots \\ \mathbf{X}_n^\top \end{pmatrix}.$$

We assume $r(\mathbb{X}) = p$. We sometimes use the notation $\boldsymbol{\beta}_0 = (\beta_{01}, \dots, \beta_{0p})^\top$ to denote the true regression parameter (but the notation $\boldsymbol{\beta}$ can also mean the true parameter).

Note. The (conditional) means of Y_1, \dots, Y_n vary because the canonical parameters $\theta_1, \dots, \theta_n$ depend on \mathbf{X}_i . The dispersion parameter φ is the same for all observations, it must not depend on \mathbf{X}_i (recall homoscedasticity in linear regression). However, the variances of Y_1, \dots, Y_n depend on the mean through the variance function $V(\mu_i)$, and hence vary with \mathbf{X}_i .

Note. The link function postulates a possibly non-linear relationship between the expectation of the response μ_i and the linear predictor $\eta_i = \mathbf{X}_i^\top \boldsymbol{\beta}$. It has to be specified in advance. There are methods to verify the choice of the link function for a specific data set (see Chapter XX). It is enough to specify the link function up to a non-zero proportionality constant (if $c \neq 0$, g and cg lead to the same model).

Definition 2.4. The link function g is called *the canonical link*[§] for the distribution f if it equates the linear predictor η_i with the canonical parameter θ_i .

* Český zobecněný lineární model † Český lineární prediktor ‡ Český linková funkce § Český kanonický link

Lemma 2.2. (Properties of canonical link)

- (i) The canonical link is equal to the inverse of b' , that is $g(\mu_i) = (b')^{-1}(\mu_i)$.
- (ii) The canonical link satisfies the equation $g'(\mu_i) = 1/V(\mu_i)$.

Note. For each distribution f from the exponential family, there is a unique (up to a non-zero proportionality constant) and specific canonical link function. Canonical link functions have certain nice properties that will become apparent later on.

Examples

Normal distribution $Y_i \sim N(\mu_i, \sigma^2)$

$$\theta_i = \mu_i, \quad \varphi = \sigma^2, \quad b(\theta_i) = \frac{\theta_i^2}{2}, \quad \mu_i = b'(\theta_i) = \theta_i, \quad \text{var } Y_i = \sigma^2, \quad V(\mu) = 1.$$

The canonical link is $g(\mu_i) = (b')^{-1}(\mu_i) = \mu_i$ (identity link).

This implies $E Y_i = \eta_i = \mathbf{X}_i^T \boldsymbol{\beta}$. This is the normal linear model.

Gamma distribution $Y_i \sim \Gamma(a_i, p)$

$$\theta_i = -\frac{a_i}{p}, \quad \varphi = 1/p, \quad b(\theta_i) = -\log(-\theta_i), \quad \mu_i = b'(\theta_i) = -1/\theta_i, \quad \text{var } Y_i = \varphi \mu_i^2.$$

The canonical link is $g(\mu_i) = (b')^{-1}(\mu_i) \propto 1/\mu_i$ (inverse link*).

This implies $E Y_i = g^{-1}(\eta_i) = 1/\mathbf{X}_i^T \boldsymbol{\beta}$.

Inverse Gaussian distribution $Y_i \sim \text{IG}(\mu_i, \lambda)$

$$\theta_i = -\frac{1}{2\mu_i^2}, \quad \varphi = 1/\lambda, \quad b(\theta_i) = -\sqrt{-2\theta_i}, \quad \mu_i = b'(\theta_i) = 1/\sqrt{-2\theta_i}, \quad \text{var } Y_i = \varphi \mu_i^3.$$

The canonical link is[†] $g(\mu_i) = (b')^{-1}(\mu_i) \propto 1/\mu_i^2$.

This implies $E Y_i = g^{-1}(\eta_i) = 1/\sqrt{\mathbf{X}_i^T \boldsymbol{\beta}}$.

Poisson distribution $Y_i \sim \text{Po}(\lambda_i)$

$$\theta_i = \log \lambda_i, \quad b(\theta_i) = \exp(\theta_i), \quad \mu_i = b'(\theta_i) = \exp(\theta_i), \quad \text{var } Y_i = \mu_i.$$

The canonical link is $g(\mu_i) = (b')^{-1}(\mu_i) = \log \mu_i$.

This implies $E Y_i = g^{-1}(\eta_i) = \exp\{\mathbf{X}_i^T \boldsymbol{\beta}\}$. This is the loglinear model.

* we drop the minus sign † we drop the constant -2

Alternative distribution $Y_i \sim \text{Alt}(p_i)$

$$\theta_i = \log \frac{p_i}{1 - p_i}, \quad b(\theta_i) = \log(1 + \exp\{\theta_i\}), \quad \mu_i = b'(\theta_i) = \frac{e^{\theta_i}}{1 + e^{\theta_i}}, \quad \text{var } Y_i = \mu_i(1 - \mu_i).$$

The canonical link is $g(\mu_i) = \log \frac{\mu_i}{1 - \mu_i}$ (the logistic link).

This implies $E Y_i = g^{-1}(\eta_i) = \frac{\exp\{\mathbf{X}_i^\top \boldsymbol{\beta}\}}{1 + \exp\{\mathbf{X}_i^\top \boldsymbol{\beta}\}}$. This is the logistic regression model.

Parametrizations of the GLM

The primary parameteres in the GLM are the regression coefficients $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top$. However, we are also interested in parametrizing the distributions of the individual Y_i 's that depend on both the primary parameters $\boldsymbol{\beta}$ and the covariates \mathbf{X}_i . This can be done in three ways:

- by the *linear predictors* η_1, \dots, η_n ;
- by the *means* $\mu_1 \equiv E Y_1, \dots, \mu_n \equiv E Y_n$;
- by the *canonical parameters* $\theta_1, \dots, \theta_n$.

The parametrizations are related to each other as follows:

- $\eta_i = \mathbf{X}_i^\top \boldsymbol{\beta}$;
- $\eta_i = g(\mu_i)$, $\mu_i = g^{-1}(\eta_i)$;
- $\mu_i = b'(\theta_i)$, $\theta_i = (b')^{-1}(\mu_i)$;
- $\eta_i = g(b'(\theta_i))$, $\theta_i = (b')^{-1}(g^{-1}(\eta_i))$; if the link g is canonical then $\eta_i = \theta_i$.

The likelihood function

Let the true dispersion parameter φ_0 be known. The likelihood function for $\boldsymbol{\beta}$ has the form

$$L(\boldsymbol{\beta} \mid \mathbf{Y}) = \prod_{i=1}^n \exp \left\{ \frac{Y_i \theta_i - b(\theta_i)}{\varphi_0} + c(Y_i, \varphi_0) \right\},$$

where $\theta_i = (b')^{-1}(g^{-1}(\mathbf{X}_i^\top \boldsymbol{\beta}))$.

The log-likelihood is

$$\ell(\boldsymbol{\beta} \mid \mathbf{Y}) = \sum_{i=1}^n \left[\frac{Y_i \theta_i - b(\theta_i)}{\varphi_0} + c(Y_i, \varphi_0) \right]. \quad (2.3)$$

The saturated model

Suppose at least one covariate is continuous and consider a model which has the largest possible number of parameters* $p = n$. This is called *the saturated model*[†]. In the saturated model, each Y_i gets its own canonical parameter θ_i , which is unrelated to the canonical parameters of the other observations. Maximizing $L(\boldsymbol{\beta} | \mathbf{Y})$ w.r.t all $\boldsymbol{\beta} \in \mathbb{R}^n$ is the same as maximizing $L(\boldsymbol{\theta} | \mathbf{Y})$ w.r.t all $\boldsymbol{\theta} \in \mathbb{R}^n$. To obtain the MLE in the saturated model, we differentiate (2.3) w.r.t. each θ_i separately and we get n equations

$$\varphi_0^{-1}[Y_i - b'(\theta_i)] = 0, \quad i = 1, \dots, n.$$

The MLE of μ_i under the saturated model is

$$\hat{\mu}_i = Y_i.$$

The fitted values $\hat{\mu}_i \equiv \hat{Y}_i$ are equal to the observed values Y_i . This model provides a “perfect fit”. However, a “perfect fit” of this kind is rarely useful.

The saturated model with $p = n$ does not satisfy the regularity assumptions of the MLE theory (the number of parameters must be constant for the theory to apply; here $p \rightarrow \infty$ as $n \rightarrow \infty$). The estimates obtained from this model are not even consistent.

The null model

The null model[‡] is the opposite extreme. It assumes $p = 1$ and $\mathbf{X}_i = 1$ so that the model includes only the intercept and all Y_i are equally distributed.

The MLE of the common canonical parameter θ of the null model is derived in Section 2.1.2. Using $\beta_0 = \eta = g(b'(\theta))$, we get the MLE of β_0 as $\hat{\beta}_n = g(b'(\hat{\theta}_n)) = g(\bar{Y}_n)$. From the central limit theorem for iid random variables and the delta method,

$$\sqrt{n}(\hat{\beta}_n - \beta_0) \xrightarrow{D} \mathbf{N}(0, \varphi_0 V(\mu_0)[g'(\mu_0)]^2),$$

where $\mu_0 = \mathbb{E} Y_i$ (compare this with (2.2)).

Neither the null model nor the saturated model are particularly interesting. We aim to build a model which has more structure than the null model, fewer parameters than the saturated model, and fits the observed data well.

2.3 Maximum Likelihood Estimation in the GLM

Let (Y_i, \mathbf{X}_i) , $i = 1, \dots, n$ be iid random vectors of dimension $p + 1$. Let $h_i(\mathbf{x})$ be the marginal density of \mathbf{X}_i (with no assumptions about it except finite second moments).

* When all covariates are discrete, the largest possible number of parameters is p^* — the number of possible distinct values of the covariate vector \mathbf{X}_i . † *Česky saturovaný model* ‡ *Česky nulový model*

Let (Y_i, \mathbf{X}_i) , $i = 1, \dots, n$, satisfy the generalized linear model (Definition 2.3) with true parameters $\boldsymbol{\beta}_0$ and φ_0 . Consider φ_0 known. Write the conditional density of Y given $\mathbf{X} = \mathbf{x}$ as $f(y | \mathbf{x}, \boldsymbol{\beta}_0, \varphi_0)$. Then the joint density of (Y_i, \mathbf{X}_i) is $f(y | \mathbf{x}, \boldsymbol{\beta}_0, \varphi_0)h_i(\mathbf{x})$, the full likelihood is

$$L^*(\boldsymbol{\beta}) = \prod_{i=1}^n f(Y_i | \mathbf{X}_i, \boldsymbol{\beta}, \varphi_0) h_i(\mathbf{X}_i)$$

and the full log-likelihood is

$$\ell^*(\boldsymbol{\beta}) = \sum_{i=1}^n \log f(Y_i | \mathbf{X}_i, \boldsymbol{\beta}, \varphi_0) + \sum_{i=1}^n \log h_i(\mathbf{X}_i).$$

Since the rightmost sum does not depend on $\boldsymbol{\beta}$, it suffices to maximize

$$\ell(\boldsymbol{\beta}) = \sum_{i=1}^n \log f(Y_i | \mathbf{X}_i, \boldsymbol{\beta}, \varphi_0). \quad (2.4)$$

This is the log-likelihood shown previously in (2.3) (without the detailed derivation and justification needed for the validity of asymptotic results).

When the covariates are random, it is not necessary to consider, know or estimate their distribution. If the covariates were constants, the log-likelihood and the score statistic would be sums of nonidentically distributed terms. Feller-Lindeberg or Lyapunov central limit theorems would have to be applied to validate the asymptotic results, and additional assumptions would have to be imposed on the covariates. The asymptotic results for constant covariates would then turn out to be the same as the results for iid data.

The core term in the log-likelihood (2.4) that we are going to maximize can be written as

$$\sum_{i=1}^n \frac{Y_i \theta_i - b(\theta_i)}{\varphi_0}, \quad (2.5)$$

where $g(\mu_i) = \mathbf{X}_i^\top \boldsymbol{\beta}$ and $\mu_i = b'(\theta_i)$. The following theorem summarizes the main results for maximum likelihood estimation of $\boldsymbol{\beta}$.

Theorem 2.3. (likelihood equations in the GLM; [Nelder and Wedderburn 1972](#)) *Let the definition of the GLM hold. Denote by $\boldsymbol{\beta}_0$ the true parameter. Let*

$$w(\mu_i) = \frac{1}{V(\mu_i)[g'(\mu_i)]^2} > 0. \quad (2.6)$$

(i) *The score function for $\boldsymbol{\beta}$ is*

$$\mathbf{U}(\boldsymbol{\beta} | Y_i) = \varphi_0^{-1} w(\mu_i) g'(\mu_i) (Y_i - \mu_i) \mathbf{X}_i,$$

where $\mu_i = g^{-1}(\mathbf{X}_i^\top \boldsymbol{\beta})$. It satisfies $\mathbf{E} \mathbf{U}(\boldsymbol{\beta}_0 | Y_i) = \mathbf{0}$.

(ii) The score statistic for $\boldsymbol{\beta}$ is

$$\mathbf{U}_n(\boldsymbol{\beta} \mid \mathbf{Y}) = \frac{1}{\varphi_0} \sum_{i=1}^n w(\mu_i) g'(\mu_i) (Y_i - \mu_i) \mathbf{X}_i.$$

(iii) The maximum likelihood estimator $\hat{\boldsymbol{\beta}}_n$ solves the system of equations

$$\sum_{i=1}^n w(\hat{\mu}_i) g'(\hat{\mu}_i) (Y_i - \hat{\mu}_i) \mathbf{X}_i = \mathbf{0}, \quad (2.7)$$

where $\hat{\mu}_i = g^{-1}(\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_n)$.

(iv) When the link g is canonical then

$$w(\mu_i) = V(\mu_i) = \frac{1}{g'(\mu_i)},$$

the score statistic can be written as

$$\mathbf{U}_n(\boldsymbol{\beta} \mid \mathbf{Y}) = \frac{1}{\varphi_0} \sum_{i=1}^n (Y_i - \mu_i) \mathbf{X}_i,$$

and the likelihood equations are

$$\sum_{i=1}^n Y_i \mathbf{X}_i = \sum_{i=1}^n \hat{\mu}_i \mathbf{X}_i.$$

Note. When the link g is canonical then $\mathbf{S} = \sum_{i=1}^n Y_i \mathbf{X}_i$ is the sufficient statistic and the MLE equates the observed value of \mathbf{S} to its expectation under the model (conditional on the covariates).

Definition 2.5. $\hat{\mu}_i = g^{-1}(\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_n)$ are called *the fitted values**.

The next step is to investigate the observed and expected information matrices for $\boldsymbol{\beta}$. Let $\mathbf{a}^{\otimes 2} \stackrel{\text{df}}{=} \mathbf{a} \mathbf{a}^\top$.

Theorem 2.4. (on information matrices in the GLM) *Let the definition of the GLM hold. Let $\mathbf{E}_{\mathbf{X}} w(\mu_i) \mathbf{X}_i^{\otimes 2}$ be finite and of full rank.*

(i) The contribution of the i -th observation to the observed information matrix is

$$I(\boldsymbol{\beta} \mid Y_i) = \frac{1}{\varphi_0} [w(\mu_i) \mathbf{X}_i^{\otimes 2} - \mathbb{J}_i],$$

where

$$\mathbb{J}_i = \left[w'(\mu_i) + w(\mu_i) \frac{g''(\mu_i)}{g'(\mu_i)} \right] (Y_i - \mu_i) \mathbf{X}_i^{\otimes 2}.$$

The observed information matrix is $I_n(\boldsymbol{\beta} \mid \mathbf{Y}) = n^{-1} \sum_{i=1}^n I(\boldsymbol{\beta} \mid Y_i)$.

* Český vyrovnané hodnoty

(ii) When evaluated at the true β_0 , $\mathbb{E} \mathbb{J}_i = 0$. The Fisher (expected) information matrix at the true β_0 is

$$I(\beta_0) = \mathbb{E} I(\beta_0 | Y_i) = \frac{1}{\varphi_0} \mathbb{E}_{\mathbf{X}} w(\mu_i) \mathbf{X}_i^{\otimes 2}. \quad (2.8)$$

By assumptions, it is finite and of full rank. It holds that $\text{var } \mathbf{U}(\beta_0 | Y_i) = I(\beta_0)$.

(iii) When the link g is canonical then $\mathbb{J}_i = 0$ at any β for all i , the observed information matrix is positive definite at all β , the log-likelihood is concave, the likelihood equations have just one solution and it is the MLE.

Note. If the link g is not canonical, there is no guarantee that a solution to the likelihood equations is the MLE. The likelihood is not concave, the equations may have multiple solutions. Numerical algorithms for solving the likelihood equations may iterate slowly and converge to the wrong solution.

The Fisher information matrix $I(\beta_0)$ can be consistently estimated by the empirical estimator

$$\hat{I}_n = \frac{1}{n\varphi_0} \sum_{i=1}^n w(\hat{\mu}_i) \mathbf{X}_i^{\otimes 2} = \frac{1}{n\varphi_0} \mathbf{X}^T \hat{\mathbb{W}} \mathbf{X}, \quad (2.9)$$

where $\hat{\mathbb{W}}$ is the $n \times n$ diagonal matrix $\text{diag}(w(\hat{\mu}_1), \dots, w(\hat{\mu}_n))$. When φ_0 is unknown it is replaced by a consistent estimator $\hat{\varphi}_n$, which will be introduced in Section 2.5.

2.4 Algorithm for Fitting the GLM

The parameters of the GLM can be estimated by a numerical algorithm called *iterative weighted least squares** [IWLS].

Theorem 2.5. (Nelder and Wedderburn 1972) The MLE $\hat{\beta}_n$ in the GLM solves the system of equations

$$\hat{\beta}_n = (\mathbf{X}^T \hat{\mathbb{W}} \mathbf{X})^{-1} (\mathbf{X}^T \hat{\mathbb{W}} \hat{\mathbf{Z}}),$$

where $\hat{\mathbb{W}} = \text{diag}(w(\hat{\mu}_1), \dots, w(\hat{\mu}_n))$, $\hat{\mathbf{Z}}$ is an n -vector with components

$$\hat{Z}_i = \hat{\eta}_i + (Y_i - \hat{\mu}_i) g'(\hat{\mu}_i),$$

$\hat{\mu}_i = g^{-1}(\hat{\eta}_i)$, and $\hat{\eta}_i = \mathbf{X}_i^T \hat{\beta}_n$.

Note. $\hat{\mathbf{Z}}$ is called *the adjusted dependent variable*[†]. Notice that \hat{Z}_i is the linear approximation to $g(Y_i)$ by Taylor expansion around $\hat{\mu}_i$. Unlike $g(Y_i)$, it can be calculated even if Y_i is outside of the domain of g .

* Český iterativní vážené nejmenší čtverce † Český upravená odezva

Note. When the link g is canonical then $\hat{\mathbb{W}} = \text{diag}(V(\hat{\mu}_1), \dots, V(\hat{\mu}_n))$ and

$$\hat{Z}_i = \hat{\eta}_i + \frac{Y_i - \hat{\mu}_i}{V(\hat{\mu}_i)}.$$

One cannot calculate $\hat{\beta}_n$ directly from Theorem 2.5 because it appears on both the left-hand side as well as the right-hand side. However, the result motivates the following iterative algorithm.

Iterative weighted least squares algorithm

Step 1. Take initial values $\hat{\mu}_i^{(0)} = Y_i$ (or $Y_i \pm \varepsilon$ if Y_i is not within the domain of g). Set $k := 0$.

Step 2. Calculate $\hat{\mathbb{W}}^{(k)} = \text{diag}(w(\hat{\mu}_1^{(k)}), \dots, w(\hat{\mu}_n^{(k)}))$ and $\hat{Z}_i^{(k)} = g(\hat{\mu}_i^{(k)}) + (Y_i - \hat{\mu}_i^{(k)})g'(\hat{\mu}_i^{(k)})$.

Step 3. Take

$$\hat{\beta}_n^{(k+1)} = (\mathbb{X}^T \hat{\mathbb{W}}^{(k)} \mathbb{X})^{-1} (\mathbb{X}^T \hat{\mathbb{W}}^{(k)} \hat{Z}^{(k)}).$$

Step 4. Calculate $\hat{\mu}_i^{(k+1)} = g^{-1}(\mathbb{X}_i^T \hat{\beta}_n^{(k+1)})$.

Step 5. Set $k := k + 1$.

Iterate steps 2–5 until convergence, for example until $\|\hat{\beta}_n^{(k)} - \hat{\beta}_n^{(k-1)}\| < \delta$, where δ is a pre-specified tolerance parameter. If the model is well formulated, the algorithm usually converges in 5–7 steps.

Note.

- The IWLS algorithm is a special case of the Fisher scoring algorithm.
- According to (2.9), the matrix $(\mathbb{X}^T \hat{\mathbb{W}}^{(k)} \mathbb{X})^{-1}$ estimates (up to a proportionality constant) the inverse information matrix. Thus, an estimate of the asymptotic variance of $\hat{\beta}_n$ is obtained by the IWLS as well (just make sure to update it after the last iteration of $\hat{\beta}_n^{(k)}$).
- Let $\mathbb{X}^* = \hat{\mathbb{W}}^{1/2} \mathbb{X}$ and $\mathbf{Y}^* = \hat{\mathbb{W}}^{1/2} \hat{Z}$. Then $\hat{\beta}_n$ can be written as an ordinary least squares estimator $\hat{\beta}_n = (\mathbb{X}^{*T} \mathbb{X}^*)^{-1} \mathbb{X}^{*T} \mathbf{Y}^*$. This is useful for extending the diagnostic methods available for the linear model to the GLM.

2.5 Estimation of the Dispersion Parameter

The dispersion parameter φ_0 is usually unknown (unless we work with Poisson or alternative distributions). This fact does not alter the estimation of β_0 or the asymptotic properties of $\hat{\beta}_n$ but we occasionally need an estimator for φ_0 . Instead of using the method of maximum likelihood, φ_0 is estimated by a modified method of moments.

Definition 2.6. The statistic

$$X^2 = \sum_{i=1}^n \frac{(Y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)} \quad (2.10)$$

is called *the Pearson chi-square statistic*^{*}. An estimator for φ_0 is given by

$$\hat{\varphi}_n = \frac{X^2}{n - p}.$$

Note. When the distribution of Y_i is normal, X^2 is the residual sum of squares SS_e and $\hat{\varphi}_n$ is the usual estimator of residual variance.

The next theorem provides conditions for consistency of $\hat{\varphi}_n$.

Theorem 2.6. Let $h(y, \mathbf{x}, \boldsymbol{\beta}) = \frac{[y - g^{-1}(\mathbf{x}^\top \boldsymbol{\beta})]^2}{V(g^{-1}(\mathbf{x}^\top \boldsymbol{\beta}))}$. Suppose there exists a function $C(y, \mathbf{x})$ such that $|\partial h / \partial \boldsymbol{\beta}| \leq C(y, \mathbf{x})$ in a neighborhood \mathcal{B}_0 of $\boldsymbol{\beta}_0$ and $\mathbb{E}C(Y_i, \mathbf{X}_i)$ exists and is finite. Then $\hat{\varphi}_n \xrightarrow{P} \varphi_0$.

Note. The condition of Theorem 2.6 is fulfilled when V and g' are bounded away from zero and V has a bounded derivative in a neighborhood of $\boldsymbol{\beta}_0$.

2.6 Deviance

Definition 2.7. The statistic

$$D(\mathbf{Y}, \hat{\boldsymbol{\beta}}_n) = 2\varphi_0 [\tilde{\ell}_n(\mathbf{Y}) - \ell_n(\hat{\boldsymbol{\beta}}_n | \mathbf{Y})],$$

where $\tilde{\ell}_n(\mathbf{Y})$ is the maximized log-likelihood of the saturated model, is called *the (un-scaled) deviance* of the model with parameters $\boldsymbol{\beta}_0 \in \mathbb{R}^p$ and observations \mathbf{Y} .

Note. In the saturated model, the MLE of μ_i is Y_i (see p. 19) and the MLE of θ_i is $\tilde{\theta}_i = (b')^{-1}(Y_i)$. The maximized log likelihood (2.5) of the saturated model is

$$\tilde{\ell}_n(\mathbf{Y}) = \frac{1}{\varphi_0} \sum_{i=1}^n [Y_i \tilde{\theta}_i - b(\tilde{\theta}_i)].$$

In the model with parameters $\boldsymbol{\beta}_0 \in \mathbb{R}^p$, the maximized log likelihood (2.5) is

$$\ell_n(\hat{\boldsymbol{\beta}}_n | \mathbf{Y}) = \frac{1}{\varphi_0} \sum_{i=1}^n [Y_i \hat{\theta}_i - b(\hat{\theta}_i)],$$

where $\hat{\theta}_i = (b')^{-1}(\hat{\mu}_i)$. Obviously, $\tilde{\ell}_n(\mathbf{Y}) \geq \ell_n(\hat{\boldsymbol{\beta}}_n | \mathbf{Y})$.

^{*} Český Pearsonovo chí kvadrát

The unscaled deviance can be expressed as

$$D(\mathbf{Y}, \hat{\boldsymbol{\beta}}_n) = 2 \sum_{i=1}^n [Y_i(\tilde{\theta}_i - \hat{\theta}_i) - b(\tilde{\theta}_i) + b(\hat{\theta}_i)]. \quad (2.11)$$

The deviance is always non-negative, does not depend on φ_0 , and is zero if and only if the model provides a “perfect fit”.

Note.

- The deviance is a goodness-of-fit measure. When the data are normal, the deviance is equal to the residual sums of squares. It generalizes the term residual sums of squares to the GLM*.
- $D^*(\mathbf{Y}, \hat{\boldsymbol{\beta}}_n, \varphi_0) = \varphi_0^{-1} D(\mathbf{Y}, \hat{\boldsymbol{\beta}}_n)$ is called *the scaled deviance*. If φ_0 is unknown, use the moment estimator $\hat{\varphi}_n$.

2.7 Asymptotic Results

Asymptotic results for the GLM follow from the general theory of maximum likelihood estimation. The theory is reviewed in “*Summary of maximum likelihood estimation theory*” (here referred to as *MLE Summary*) available from the course website.

The following theorem transcribes the results of Theorems 2–5 from *MLE Summary* in the context of the GLM. The regularity conditions R1–R4 are assured by the specification of the model. Condition R6 has been verified in Theorem 2.3, part (i) and Theorem 2.4, part (ii).

The Fisher information matrix

$$I(\boldsymbol{\beta}_0) = \mathbf{E} I(\boldsymbol{\beta}_0 | Y_i) = \text{var} \mathbf{U}(\boldsymbol{\beta}_0 | Y_i) = \frac{1}{\varphi_0} \mathbf{E}_{\mathbf{X}} w(\mu_i) \mathbf{X}_i^{\otimes 2}$$

is finite and of full rank by assumptions imposed on the covariates (finiteness of all necessary moments and linear independence of covariates).

Theorem 2.7.

(i) The MLE $\hat{\boldsymbol{\beta}}_n$ is consistent (as long as the likelihood equations (2.7) have a unique solution).

(ii)

$$\frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\beta}_0) \xrightarrow{D} \mathbf{N}_p(\mathbf{0}, I(\boldsymbol{\beta}_0)).$$

(iii)

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) \xrightarrow{D} \mathbf{N}_p(\mathbf{0}, I^{-1}(\boldsymbol{\beta}_0)).$$

* The Pearson X^2 is another generalization.

(iv)

$$2 \log \frac{L_n(\hat{\boldsymbol{\beta}}_n | \mathbf{Y})}{L_n(\boldsymbol{\beta}_0 | \mathbf{Y})} \xrightarrow{D} \chi_p^2.$$

The information matrix $I(\boldsymbol{\beta}_0)$ can be consistently estimated by

$$\hat{I}_n = \frac{1}{n\hat{\varphi}_n} \mathbf{X}^\top \hat{\mathbb{W}} \mathbf{X}.$$

According to part (iii) of Theorem 2.7, the estimated asymptotic variance of $\hat{\boldsymbol{\beta}}_n$ is

$$\hat{I}_n^{-1}/n = \hat{\varphi}_n (\mathbf{X}^\top \hat{\mathbb{W}} \mathbf{X})^{-1}. \quad (2.12)$$

Denote $\hat{\Sigma} \equiv (\mathbf{X}^\top \hat{\mathbb{W}} \mathbf{X})^{-1}$ so that $\hat{\varphi}_n \hat{\Sigma}$ estimates $\text{var } \hat{\boldsymbol{\beta}}_n$.

Let us consider the problem of testing the simple hypothesis

$$H_0 : \boldsymbol{\beta} = \boldsymbol{\beta}_0 \quad \text{against} \quad H_1 : \boldsymbol{\beta} \neq \boldsymbol{\beta}_0.$$

The test statistics and their null distributions are established by the following theorem, which is based on Definition 5 and Theorem 7 from *MLE Summary*.

Theorem 2.8.

(i) **Score (Rao) test.** Let $\mu_i^0 = g^{-1}(\mathbf{X}_i^\top \boldsymbol{\beta}_0)$, $\mathbb{W}^0 = \text{diag}(w(\mu_1^0), \dots, w(\mu_n^0))$, denote $\Sigma^0 = (\mathbf{X}^\top \mathbb{W}^0 \mathbf{X})^{-1}$. If H_0 holds then

$$\begin{aligned} R_n &= \frac{1}{n} \mathbf{U}_n(\boldsymbol{\beta}_0)^\top \hat{I}_n^{-1} \mathbf{U}_n(\boldsymbol{\beta}_0) \\ &= \frac{1}{\hat{\varphi}_n} \left(\sum_{i=1}^n w(\mu_i^0) g'(\mu_i^0) (Y_i - \mu_i^0) \mathbf{X}_i \right)^\top \Sigma^0 \left(\sum_{i=1}^n w(\mu_i^0) g'(\mu_i^0) (Y_i - \mu_i^0) \mathbf{X}_i \right) \\ &\xrightarrow{D} \chi_p^2 \end{aligned}$$

(ii) **Wald test.** If H_0 holds then

$$W_n = n(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0)^\top \hat{I}_n (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) = \frac{1}{\hat{\varphi}_n} (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0)^\top \hat{\Sigma}^{-1} (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) \xrightarrow{D} \chi_p^2$$

(iii) **Likelihood ratio test.** Let $\theta_i^0 = (b')^{-1}(\mu_i^0)$. If H_0 holds then

$$\lambda_n = 2[\ell_n(\hat{\boldsymbol{\beta}}_n | \mathbf{Y}) - \ell_n(\boldsymbol{\beta}_0 | \mathbf{Y})] = \frac{2}{\hat{\varphi}_n} \sum_{i=1}^n [Y_i(\hat{\theta}_i - \theta_i^0) - b(\hat{\theta}_i) + b(\theta_i^0)] \xrightarrow{D} \chi_p^2$$

The simple hypothesis is rarely of interest for applications. We are more interested in composite hypotheses, for example, in testing that the last m components of the

regression parameter vector are all zero (without loss of generality: the components of β can be always rearranged in this way). Take

$$H_0^* : \begin{pmatrix} \beta_{p-m+1} \\ \beta_{p-m+2} \\ \vdots \\ \beta_p \end{pmatrix} = \mathbf{0} \quad \text{against} \quad H_1^* : \begin{pmatrix} \beta_{p-m+1} \\ \beta_{p-m+2} \\ \vdots \\ \beta_p \end{pmatrix} \neq \mathbf{0}$$

for some $m < p$. If H_0^* is true then the last m parameters attain zero value and the last m columns of the covariate vector can be excluded from the model. The null hypothesis specifies a submodel (with $p - m$ parameters) of the full model with (p parameters).

Denote $\beta_M = (\beta_{p-m+1}, \dots, \beta_p)^\top$ and $\mathbf{X}_i^M = (X_{p-m+1}, \dots, X_p)^\top$. Let $\hat{\beta}_M = (\hat{\beta}_{p-m+1}, \dots, \hat{\beta}_p)^\top$ be the MLE of β_M under the larger model. Let $\tilde{\beta}_n$ be the MLE of β under the submodel (subject to the constraint $\beta_M = \mathbf{0}$), let $\tilde{\mu}_i = g^{-1}(\mathbf{X}_i^\top \tilde{\beta}_n)$ be the fitted values under the submodel.

Partition the $p \times p$ matrix $\hat{\Sigma} = \hat{I}_n^{-1}/(n\hat{\varphi}_n) = (\mathbb{X}^\top \hat{\mathbb{W}} \mathbb{X})^{-1}$ (the estimated asymptotic variance of $\hat{\beta}_n$ without $\hat{\varphi}_n$) into four blocks

$$\hat{\Sigma} = \begin{pmatrix} \hat{\Sigma}_A & \hat{\Sigma}_B \\ \hat{\Sigma}_B^\top & \hat{\Sigma}_M \end{pmatrix},$$

where the lower right block $\hat{\Sigma}_M$ is of size $m \times m$.

Theorem 2.9.

- (i) **Score (Rao) test.** Let $\tilde{\mathbb{W}} = \text{diag}(w(\tilde{\mu}_1), \dots, w(\tilde{\mu}_n))$. Let $\tilde{\Sigma}_M$ be the $m \times m$ lower right block of the matrix $\tilde{\Sigma} = (\mathbb{X}^\top \tilde{\mathbb{W}} \mathbb{X})^{-1}$. Denote by $\tilde{\varphi}_n$ the estimator of the dispersion parameter calculated under the submodel (under H_0^*). If H_0^* holds then

$$R_n^* = \frac{1}{\tilde{\varphi}_n} \left(\sum_{i=1}^n w(\tilde{\mu}_i) g'(\tilde{\mu}_i) (Y_i - \tilde{\mu}_i) \mathbf{X}_i^M \right)^\top \tilde{\Sigma}_M \left(\sum_{i=1}^n w(\tilde{\mu}_i) g'(\tilde{\mu}_i) (Y_i - \tilde{\mu}_i) \mathbf{X}_i^M \right) \xrightarrow{D} \chi_m^2.$$

- (ii) **Wald test.** Denote by $\hat{\varphi}_n$ the estimator of the dispersion parameter calculated under the larger model (not assuming that H_0^* is true). If H_0^* holds then

$$W_n^* = \frac{1}{\hat{\varphi}_n} (\hat{\beta}^M)^\top \hat{\Sigma}_M^{-1} (\hat{\beta}^M) \xrightarrow{D} \chi_m^2.$$

- (iii) **Likelihood ratio (deviance) test.** Let $D(\mathbf{Y} \mid \tilde{\beta})$ be the (unscaled) deviance of the submodel, let $D(\mathbf{Y} \mid \hat{\beta})$ be the (unscaled) deviance of the larger model. Let the estimate $\hat{\varphi}_n$ be calculated under the larger model (not assuming that H_0^* is true). If H_0^* holds then

$$\lambda_n^* = \frac{1}{\hat{\varphi}_n} [D(\mathbf{Y} \mid \tilde{\beta}) - D(\mathbf{Y} \mid \hat{\beta})] \xrightarrow{D} \chi_m^2.$$

Note.

- Theorem 2.9 follows from Definition 6 and Theorem 9 in *MLE Summary*. The hypothesis H_0^* is rejected at the asymptotic level of α if the chosen test statistic (it must be selected in advance) exceeds the $1 - \alpha$ quantile of the χ_m^2 distribution.
- Under the standard linear regression model with normal distribution, these three test statistics are all equal to the F test statistic (1.1) for submodel testing. In that case, the exact distribution of the test statistics under the null hypothesis is $F_{m,n-p}$. When normality does not hold or the link is not identity, the three test statistics are not the same and we only know that their asymptotic distribution is χ_m^2 .
- Generally, the likelihood ratio test statistic is twice the difference in the log likelihoods between the model and the submodel. However, it can be also expressed as a properly scaled difference in deviances between the submodel and the model. *The deviance test is the preferred tool for testing submodels in generalized linear models.*
- The Wald and Rao statistics are asymptotically equivalent to the likelihood ratio test statistic. However, in finite samples they may be different. Unlike the likelihood ratio test statistic, the Wald test statistic depends on the parametrization of the model and tends to have the slowest convergence to the asymptotic distribution. For these reasons, the Wald statistic is the least desirable of the three.
- An important special case is $m = 1$ (testing of a single parameter). Then the Wald statistic for testing zero value of the j -th parameter is

$$\left(\frac{\hat{\beta}_j}{\sqrt{\hat{\varphi}_n \hat{\sigma}_{jj}^2}} \right)^2, \tag{2.13}$$

where $\hat{\sigma}_{jj}^2$ is the j -th diagonal element of $\hat{\Sigma}$. Before applying the square, these statistics are asymptotically standard normal; in this form they are automatically provided in the output of almost any statistical software for fitting the GLM.

- The deviance of the current model $D(\mathbf{Y} \mid \hat{\beta})$ is twice the difference in log likelihoods between the saturated model and the current model. However, the deviance cannot be in general used as a test statistic to compare the goodness-of-fit of the current model to the saturated model unless all covariates are discrete (otherwise the number of parameters of the saturated model grows to infinity and Theorem 9 from *MLE Summary* does not hold). Differences in deviances between a submodel and a larger model do not have this problem.

Confidence intervals

The simplest confidence intervals for the individual parameters are based on Wald test statistics (2.13). The interval with end points

$$\hat{\beta}_j \pm u_{1-\alpha/2} \sqrt{\hat{\varphi}_n \hat{\sigma}_{jj}^2},$$

covers β_j with probability converging to $1 - \alpha$.

Better confidence intervals would be obtained from inverting acceptance regions of the Rao or likelihood ratio test statistics or using profile likelihood methods.

Wald-type confidence intervals for linear combinations of parameters $\mathbf{c}^\top \boldsymbol{\beta}_0$ where $\mathbf{0} \neq \mathbf{c} \in \mathbb{R}^p$ can be obtained easily from Theorem 2.7 part (iii). An asymptotic confidence interval with coverage probability converging to $1 - \alpha$ has end points

$$\mathbf{c}^\top \hat{\boldsymbol{\beta}}_n \pm u_{1-\alpha/2} \sqrt{\hat{\varphi}_n \mathbf{c}^\top \hat{\Sigma} \mathbf{c}}.$$

2.8 Diagnostic Methods for the GLM

Diagnostic methods can be derived from the linear model using Theorem 2.5. Let $\mathbb{X}^* = \hat{\mathbb{W}}^{1/2} \mathbb{X}$ and $\mathbf{Y}^* = \hat{\mathbb{W}}^{1/2} \hat{\mathbf{Z}}$. Write $\hat{\boldsymbol{\beta}}_n$ as an ordinary least squares estimator $\hat{\boldsymbol{\beta}}_n = (\mathbb{X}^{*\top} \mathbb{X}^*)^{-1} \mathbb{X}^{*\top} \mathbf{Y}^*$. Let

$$\mathbb{H}^* = \mathbb{X}^* (\mathbb{X}^{*\top} \mathbb{X}^*)^{-1} \mathbb{X}^{*\top} = \hat{\mathbb{W}}^{1/2} \mathbb{X} (\mathbb{X}^\top \hat{\mathbb{W}} \mathbb{X})^{-1} \mathbb{X}^\top \hat{\mathbb{W}}^{1/2},$$

and $\hat{\mathbf{Y}}^* = \mathbb{X}^* \hat{\boldsymbol{\beta}} = \mathbb{H}^* \mathbf{Y}^* = \hat{\mathbb{W}}^{1/2} \mathbb{X} (\mathbb{X}^\top \hat{\mathbb{W}} \mathbb{X})^{-1} \mathbb{X}^\top \hat{\mathbb{W}} \hat{\mathbf{Z}}$.

2.8.1 Pearson residuals

Pearson residuals are defined by the identity $\mathbf{Y}^* - \hat{\mathbf{Y}}^* = \hat{\mathbb{W}}^{1/2} \hat{\mathbf{Z}} - \hat{\mathbb{W}}^{1/2} \mathbb{X} \hat{\boldsymbol{\beta}}$, which gives the following residuals for the individual observations

$$r_i^P = \frac{Y_i - \hat{\mu}_i}{\sqrt{V(\hat{\mu}_i)}}.$$

Sum of squares of Pearson residuals is equal to the Pearson X^2 statistic:

$$\sum_{i=1}^n (r_i^P)^2 = X^2.$$

2.8.2 Leverages

It can be shown that

$$\text{var } r_i^P \approx \varphi_0(1 - h_{ii}^*),$$

where h_{ii}^* , the i -th diagonal element of \mathbb{H}^* , is called *the leverage*. Potentially influential observations can be identified by the rule of thumb $h_{ii}^* > 2p/(n - 2p)$. These observations are sort of atypical in their covariates and thus may have unduly strong influence on the results of the model fit.

2.8.3 Standardized Pearson residuals

Standardized Pearson residuals divide r_i^P by the square root of the estimated approximate variance of r_i^P :

$$r_i^{PS} = \frac{Y_i - \hat{\mu}_i}{\sqrt{\hat{\varphi}_n V(\hat{\mu}_i)(1 - h_{ii}^*)}}.$$

They have approximately unit variance.

2.8.4 Deviance residuals

Deviance residuals are signed square roots of the contributions of the observations to the deviance. Let $\tilde{\theta}_i = (b')^{-1}(Y_i)$, $d_i = 2\{Y_i[\tilde{\theta}_i - \hat{\theta}_i] - b(\tilde{\theta}_i) + b(\hat{\theta}_i)\}$, and define the deviance residual as

$$r_i^D = \text{sgn}(Y_i - \hat{\mu}_i)\sqrt{d_i}.$$

Sum of squares of deviance residuals is equal to the deviance:

$$\sum_{i=1}^n (r_i^D)^2 = D(\mathbf{Y} \mid \hat{\boldsymbol{\beta}}).$$

2.8.5 Standardized deviance residuals

Standardized deviance residuals use the same normalization as standardized Pearson residuals.

$$r_i^{DS} = \frac{\text{sgn}(Y_i - \hat{\mu}_i)\sqrt{d_i}}{\sqrt{\hat{\varphi}_n(1 - h_{ii}^*)}}.$$

These are the default residuals in R.

2.8.6 Cook's distance

Cook's distance measures the influence of the i -th observation on the estimates of regression parameters $\hat{\beta}$. Let $\hat{\beta}_{(i)}$ denote the estimates calculated after deletion of the i -th observation from the data set. Cook's distance is defined as

$$CD_i = \frac{1}{p\hat{\varphi}_n} (\hat{\beta} - \hat{\beta}_{(i)})^\top \mathbb{X}^{*\top} \mathbb{X}^* (\hat{\beta} - \hat{\beta}_{(i)}).$$

In linear regression, it can be shown that

$$CD_i = \frac{1}{p\hat{\varphi}_n} \left(\frac{Y_i^* - \hat{Y}_i^*}{\sqrt{1 - h_{ii}^*}} \right)^2 \frac{h_{ii}^*}{1 - h_{ii}^*} = \frac{1}{p} (r_i^{PS})^2 \frac{h_{ii}^*}{1 - h_{ii}^*}.$$

This is how Cook's distance is calculated in the GLM. An observation is considered influential if $CD_i > \frac{8}{n-2p}$.

2.8.7 Residual plots

Residual plots are created and used in a direct analogy with the linear model. However, for some data types (e.g. binary data) the residual plots are much less informative and require smoothing to yield any useful information. In general, residual plots are somewhat less useful in the GLM than they are in the linear model.

2.8.8 Diagnostics of the link function

We only mention two simple methods for checking that the correct link function was selected. Plotting the adjusted dependent variable \hat{Z}_i against the linear predictor $\hat{\eta}_i$ provides a graphical check. If the link is correct the plot should reveal a linear pattern. A formal test can be obtained by adding $(\hat{\eta}_i)^2$ to the model as an additional covariate and testing that its parameter is zero. If the hypothesis is rejected the link may be incorrect.

Both methods are sensitive to inappropriate transformations of the regressors. If the transformations are not chosen well, both methods may indicate a problem even if the link is correct.

Incorrect link functions do not have a serious effect on deciding which regressors affect the response or on the results of submodel testing. The choice of the link function is important if the primary goal of the analysis is prediction.

2.9 Model-building strategies

Model-building strategies for generalized linear models do not differ from the strategies applied to other regression models, including linear regression. The primary tool for

model building are deviance tests comparing a larger model with a submodel. If the deviance test is significant it means that the terms in the larger model cannot be removed without a significant decrease in the quality of model fit.

Since the development of the final model usually involves repeated applications of deviance tests, each performed on a selected level α (usually $\alpha = 0.05$), it is clear that the overall procedure does not preserve the desired level. If many tests are done then the final model is likely to include terms that in fact do not affect the response at all (*overfitting*). There is no universal and reliable method for adjusting the levels of the individual tests so that the overall probability of including irrelevant terms is under control. Nevertheless the analyst should be aware of this problem and should not interpret the p-values of submodel tests too dogmatically.

Approaches for developing reasonable models vary with the nature of the problem, structure of the data and questions to be addressed by the analysis. There is no universal solution to be recommended. Each problem requires careful consideration by the analyst taking into account the nature of the problem, the data-collection methods and tools, the meaning of the variables included in the dataset, their mutual relationships, and the goals of the analysis.

If *prediction* is the primary goal, it is useful to consider rich and flexible models. Omission of an important term from the model or its inclusion with an inappropriate transformation may have detrimental biasing effects on the predictions. If unnecessary covariates are left in, the variability in the predicted response is increased but the predictions are not biased. Interpretation of regression parameters is usually not that important. In prediction analyses, validation of the prediction model should be performed either by dividing the data set into disjoint training (used for model building) and validation (used for evaluation of the predictions) subsets or at least by cross-validation (predictions of each observation by a model fitted on data excluding that observation). Validation is a very useful tool for selection of the best prediction model out of several candidates.

If the goal is to *evaluate covariate effects* (“how does covariate X affect the mean of the response Y ?”), one must be really careful about several things. First, the covariate of interest must be kept in the model even if it is not significant – otherwise its effect cannot be evaluated. Second, the regression parameters expressing the influence of the covariate of interest should have a straightforward interpretation. Thus, we cannot afford to model the effect of X by a complicated function that cannot be easily summarized (splines of order > 1 , polynomials), or to use complex transformations of the response or link functions that are difficult to interpret. Third, there might be covariates that should be kept in the model regardless of their significance (suspected confounders) and/or covariates that should not be included in the model no matter how significant they are (variables on the causal pathway between X and Y , variables that are influenced by the value of Y). Thus, making reasonable decisions about which covariates should be included in the model and which should be dropped is not based solely on significance tests but also

on external expert knowledge of the problem to be analyzed. It is precisely this issue that makes automated computer-based algorithms (unsupervised stepwise regression, regression trees, neural networks, deep learning, etc.) unable to solve certain problems acceptably.

Another common problem in model-building strategies is the inclusion of *interactions*, especially when the number of covariates that can be considered for interactions is quite large. The strategy that starts with a model that includes a lot of main effects as well as all possible two-way interactions between them, and tries to gradually eliminate the superfluous terms usually does not lead to a good model. With this approach, we are likely to end up with a model that suffers from overfitting, keeps a lot of unnecessary interactions and is hard to interpret. It is better to fit only the main effects first, eliminate those that are not contributing to the model, and then try to add two-way interactions of the remaining terms one by one. This strategy is much more likely to end up only with interactions that really matter. Considering higher order interactions (three-way, four-way, . . .) is usually a hopeless task. It is better not to consider them at all, except in analyses where, for some reason, such interactions are among the terms of interest.

There is one principle about building models with interactions, which is almost universally valid and the analyst should take care not to violate it. The models should be built *hierarchically*, meaning that if a covariate is present in a higher-order interaction, then all its corresponding lower-order interactions as well as the main effects should be included in the model as well, no matter if they are significant or not. This principle should be ignored only in analyses where there is a sound justification for its violation.

This brief exposition of model-building strategies cannot be complete and should be understood in the whole context of the particular task to be done. As noted earlier, each problem should be carefully considered in order to choose a tailor-made strategy that works well for it. This requires practical experience. The analyst should be aware that there is no such thing as the true model and that his task is not to discover it. All models are wrong – we are only looking for an acceptable model that provides satisfactory answers to the questions of interest.

3 Generalized Linear Model for Discrete Responses

3.1 Analysis of Binary Data

3.1.1 Alternative vs. binomial data

Let $Y_{ij}^* \sim \text{Alt}(\pi_i)$, $\pi \in (0, 1)$, be independent variables for $i = 1, \dots, K$, $j = 1, \dots, m_i$. For a fixed i , $Y_{i1}^*, \dots, Y_{im_i}^*$ are identically distributed. The total number of observations is $N = \sum_{i=1}^K m_i$. Let π_i depend on $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^\top$ through the linear predictor $\eta_i = \mathbf{X}_i^\top \boldsymbol{\beta}$, $\boldsymbol{\beta}$ is the vector of unknown regression coefficients to be estimated. Therefore $Y_{i1}^*, \dots, Y_{im_i}^*$ share the same covariate vector \mathbf{X}_i .

The response $Y_{ij}^* \sim \text{Alt}(\pi_i)$ has a distribution of exponential family with $\mu_i \equiv \mathbb{E} Y_{ij}^* = \pi_i$ and $\text{var} Y_{ij}^* = \pi_i(1 - \pi_i)$. The variance function is $V(\mu) = \mu(1 - \mu)$, the dispersion parameter is $\varphi = 1$, the canonical parameter is $\theta_i = \log \frac{\pi_i}{1 - \pi_i}$. Finally, $b(\theta_i) = \log(1 + e^{\theta_i}) = \log \frac{1}{1 - \pi_i}$.

Denote $Y_i = \sum_{j=1}^{m_i} Y_{ij}^*$. Then $Y_i \sim \text{Bi}(m_i, \pi_i)$. Thus, binomial responses can be treated together with alternative responses even though the binomial distribution does not strictly belong to the exponential family.

The dataset for alternative/binomial regression can be arranged in two different ways. It is recommended not to mix the two data formats in a single dataset.

Format A. The dataset is arranged so that there are N rows and each value of \mathbf{X}_i appears in m_i different rows. The row corresponding to the ij -th observation includes Y_{ij}^* and \mathbf{X}_i . This is the Bernoulli format of the data.

Format B. The dataset is arranged so that there are K rows and each value of \mathbf{X}_i appears only once. The i -th row includes Y_i , m_i , and \mathbf{X}_i . This is the binomial format of the data.

There are two different kinds of asymptotics when $N \rightarrow \infty$.

1. K is constant, $m_i \rightarrow \infty$ at the same rate for all i . This happens when all covariates are discrete with a finite support.
2. $K \rightarrow \infty$, m_i are small (typically $m_i = 1$). This happens when at least one covariate is continuous.

Most of the results are the same for both kinds of asymptotics but there are certain

important differences that will be pointed out later.

3.1.2 Link functions for binary data

Because $\mu_i \equiv \pi_i \in (0, 1)$, suitable link functions are maps $(0, 1) \rightarrow \mathbb{R}$. Any quantile function of a continuous distribution on \mathbb{R} could be used as a link function for binary responses. Here are some examples:

Logistic link

Take the quantile function of the standard logistic distribution.

$$g(\mu_i) = \log \frac{\mu_i}{1 - \mu_i}, \quad \mu_i = \frac{\exp\{\mathbf{X}_i^\top \boldsymbol{\beta}\}}{1 + \exp\{\mathbf{X}_i^\top \boldsymbol{\beta}\}}.$$

This is the logistic link, the canonical link function, the most commonly used link for binary data. The model is called *the logistic regression model*^{*}.

Probit link

Take the quantile function of the standard normal distribution.

$$g(\mu_i) = \Phi^{-1}(\mu_i), \quad \mu_i = \Phi(\mathbf{X}_i^\top \boldsymbol{\beta}).$$

This is the probit link, the model is called *the probit regression model*[†]. It is used in threshold analysis, toxicology and pharmacokinetics.

Cauchit link

Take the quantile function of the standard Cauchy distribution.

$$g(\mu_i) = \tan[\pi(\mu_i - 0.5)], \quad \mu_i = \frac{1}{\pi} \arctan(\mathbf{X}_i^\top \boldsymbol{\beta}) + \frac{1}{2}.$$

This is the cauchit link, the model is called *the cauchit regression model*[‡]. It is suitable when π_i converges to 0 (1) extremely slowly for $\eta_i \rightarrow \pm\infty$.

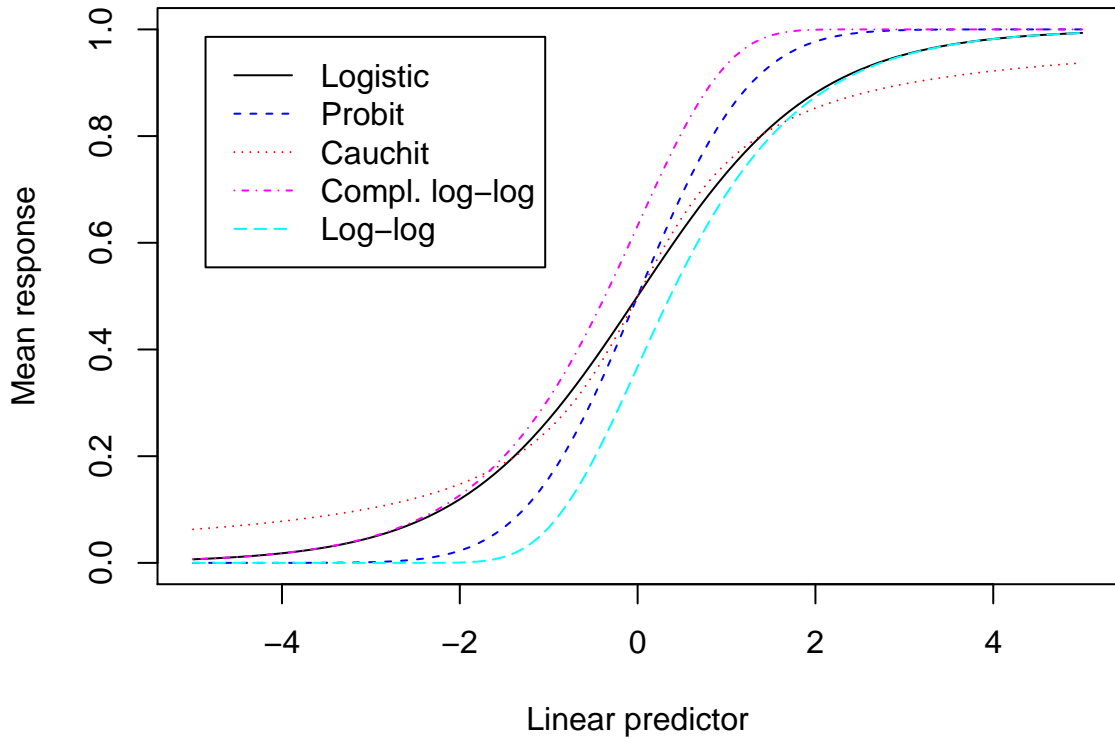
Complementary log-log link

Take the quantile function of the negative Gumbel (extreme value) random variable.

$$g(\mu_i) = \log(-\log(1 - \mu_i)), \quad \mu_i = 1 - e^{-\exp\{\mathbf{X}_i^\top \boldsymbol{\beta}\}}.$$

^{*} Český *logistická regrese* [†] Český *probitová regrese* [‡] Český *cauchitová regrese*

Figure 3.1: Inverse link functions for binary data. The linear predictor η_i is on the horizontal axis, the success probability π_i is on the vertical axis.



This link function does not possess symmetry properties. It is used in the analysis of discrete survival data. Its counterpart is the log-log link

$$g(\mu_i) = -\log(-\log(\mu_i)), \quad \mu_i = e^{-\exp\{-\mathbf{X}_i^T \boldsymbol{\beta}\}}.$$

The inverse link functions are plotted in Figure 3.1. The choice of the link function should be governed by the desired interpretation of the fitted model rather than by the data. The canonical logistic link should be the first choice unless a different interpretation is needed or there is a strong prior reason to choose a different link.

3.1.3 Binary data likelihood

There are two different sampling schemes to be considered.

- (i) Alternative responses are observed independently of each other together with the

covariates. Then m_i are random variables.

- (ii) m_i is fixed in advance, then m_i independent observations are obtained for each combination of the covariates.

The likelihoods for these two schemes only differ by a constant that does not affect the analysis. If m_i is random then the likelihood is a product of independent alternative distributions

$$\prod_{i=1}^K \prod_{j=1}^{m_i} \pi_i^{Y_{ij}^*} (1 - \pi_i)^{1 - Y_{ij}^*} = \prod_{i=1}^K \pi_i^{Y_i} (1 - \pi_i)^{m_i - Y_i}.$$

If m_i is fixed then the likelihood is a product of independent binomial distributions

$$\prod_{i=1}^K \binom{m_i}{Y_i} \pi_i^{Y_i} (1 - \pi_i)^{m_i - Y_i} = \prod_{i=1}^K \pi_i^{Y_i} (1 - \pi_i)^{m_i - Y_i} \prod_{i=1}^K \binom{m_i}{Y_i}.$$

The product of the binomial numbers does not include the parameters, so it is not relevant. The first scheme follows the framework of independent observations from a distribution of exponential type so the theory of Chapter 2 applies. The second scheme does not follow this framework but all the results have exactly the same form and properties. Therefore we do not have to distinguish the sampling schemes.

3.1.4 Threshold analysis by probit regression

The probit link has an interesting application in threshold analysis of normally distributed data.

Consider random variables U_i that follow the normal linear regression model

$$U_i = \mathbf{Z}_i^T \boldsymbol{\alpha} + \varepsilon_i, \tag{3.1}$$

where \mathbf{Z}_i are p -dimensional covariate vectors, $\boldsymbol{\alpha}$ are regression coefficients and $\varepsilon_i \sim \mathbf{N}(0, \sigma^2)$ are error terms. Now suppose that the responses U_i cannot be observed directly. Instead, a threshold C_i is provided and we learn whether the unobserved response U_i exceeds the threshold or not.

Assume that C_i is independent of U_i . The observed response is $Y_i = I(U_i < C_i)$, together with the values of the covariates \mathbf{Z}_i and the threshold C_i . The goal is to estimate the regression coefficients $\boldsymbol{\alpha}$ and the residual variance σ^2 of the underlying linear regression model (3.1).

The observations come in the form of iid triplets (Y_i, C_i, \mathbf{Z}_i) . The response Y_i follows an alternative distribution with

$$P[Y_i = 1] \equiv p_i = P[U_i < C_i].$$

Conditionally on the value of the observed threshold C_i , we get

$$p_i = P\left[\frac{U_i - \mathbf{Z}_i^T \boldsymbol{\alpha}}{\sigma} < \frac{C_i - \mathbf{Z}_i^T \boldsymbol{\alpha}}{\sigma}\right] = \Phi\left(\frac{C_i}{\sigma} - \mathbf{Z}_i^T \frac{\boldsymbol{\alpha}}{\sigma}\right).$$

Define

$$\mathbf{X}_i = \begin{pmatrix} \mathbf{Z}_i \\ C_i \end{pmatrix} \quad \text{and} \quad \boldsymbol{\beta} = \begin{pmatrix} -\boldsymbol{\alpha}/\sigma \\ 1/\sigma \end{pmatrix}.$$

This translates the problem into binary probit regression model with the linear predictor $\mathbf{X}_i^\top \boldsymbol{\beta}$. The parameters $\boldsymbol{\beta}$ can be estimated by the usual procedures for the analysis of the GLM. The parameters of interest can be obtained from $\hat{\boldsymbol{\beta}}$ as $\hat{\sigma}^2 = \frac{1}{\hat{\beta}_{p+1}^2}$ and $\hat{\alpha}_j = -\frac{\hat{\beta}_j}{\hat{\beta}_{p+1}}$.

Of course, this can only be done if the threshold values C_i are linearly independent of the covariates \mathbf{Z}_i . For example, if C_i are all set to the same value, the intercept term cannot be distinguished from the residual variance and the parameters of the original linear regression model cannot be determined.

3.1.5 Logistic regression

The logistic regression model is the most commonly used model for the analysis of binary and binomial responses.

The logistic link has the form $g(\pi_i) = \log \frac{\pi_i}{1-\pi_i}$, where $\pi_i/(1-\pi_i)$ is the odds of success. The success probabilities can be expressed as $\pi_i = \frac{\exp\{\mathbf{X}_i^\top \boldsymbol{\beta}\}}{1+\exp\{\mathbf{X}_i^\top \boldsymbol{\beta}\}}$.

Interpretation of regression parameters

Let $\mathbf{X}_i^\top \boldsymbol{\beta} = \beta_1 + \beta_2 X_2 + \dots + \beta_p X_p$. Denote $\pi_0 = P(Y_{ij}^* = 1 | X_2 = \dots = X_p = 0)$. Then

$$\log \frac{\pi_0}{1-\pi_0} = \beta_1$$

so e^{β_1} is the odds of success for an individual with zero values in all covariates.

Now consider two individuals: one with observed covariates $\mathbf{x}^0 = (1, x_2, \dots, x_p)^\top$, the other with observed covariates $\mathbf{x}^j = (1, x_2, \dots, x_j + 1, \dots, x_p)^\top$ (the j -th covariate is increased by 1, the others are the same). Denote $\pi_{X0} = P(Y_{ij}^* = 1 | \mathbf{X} = \mathbf{x}^0)$ and $\pi_{Xj} = P(Y_{ij}^* = 1 | \mathbf{X} = \mathbf{x}^j)$. Then

$$\boldsymbol{\beta}^\top \mathbf{x}^0 = \log \frac{\pi_{X0}}{1-\pi_{X0}} \quad \text{and} \quad \boldsymbol{\beta}^\top \mathbf{x}^j = \boldsymbol{\beta}^\top \mathbf{x}^0 + \beta_j = \log \frac{\pi_{Xj}}{1-\pi_{Xj}}.$$

It follows that

$$\beta_j = \log \left(\frac{\pi_{Xj}}{1-\pi_{Xj}} \cdot \frac{1-\pi_{X0}}{\pi_{X0}} \right) \quad \text{and} \quad e^{\beta_j} = \frac{\pi_{Xj}(1-\pi_{X0})}{\pi_{X0}(1-\pi_{Xj})}.$$

Thus e^{β_j} is the odds ratio for success comparing two individuals differing by one unit in the covariate X_j . E.g., if $\beta_j = 0.431$ one can say that a unit increase in the covariate

X_j increases the odds of success $e^{0.431} = 1.539$ times (or by 53.9%). When $\beta_j = 0$ the odds ratio is 1 and the covariate has no effect on the odds of success (or the probability of success) given the other covariates.

Consider a two-by-two contingency table of conditional probabilities

Covariates	$Y = 1$	$Y = 0$
$\mathbf{X} = \mathbf{x}^j$	π_{Xj}	$1 - \pi_{Xj}$
$\mathbf{X} = \mathbf{x}^0$	π_{X0}	$1 - \pi_{X0}$

The odds ratio e^{β_j} describes the association between \mathbf{X} and Y in this contingency table. The odds ratio is one if and only if there is independence.

Estimation of parameters

By Theorem 2.3, the score statistic with the canonical link is

$$\mathbf{U}_n(\boldsymbol{\beta} \mid \mathbf{Y}) = \sum_{i=1}^K \sum_{j=1}^{m_i} (Y_{ij}^* - \pi_i) \mathbf{X}_i = \sum_{i=1}^K (Y_i - m_i \pi_i) \mathbf{X}_i$$

and $\widehat{\boldsymbol{\beta}}_n$ solves the equations

$$\sum_{i=1}^K Y_i \mathbf{X}_i = \sum_{i=1}^K m_i \widehat{\pi}_i \mathbf{X}_i,$$

where

$$\widehat{\pi}_i = \frac{\exp\{\mathbf{X}_i^\top \widehat{\boldsymbol{\beta}}_n\}}{1 + \exp\{\mathbf{X}_i^\top \widehat{\boldsymbol{\beta}}_n\}}.$$

The IWLS algorithm can be implemented in two different ways depending on the data format. With the Bernoulli format A, the regression matrix \mathbb{X} includes each observed value of \mathbf{X}_i in m_i different rows, and its dimension is $N \times p$. Suppose the observations ij are ordered by the two indices $11, \dots, 1m_1, 21, \dots, 2m_2, \dots, Km_K$. Let

$$\widehat{\mathbb{W}}^{(k)} = \text{diag}(\widehat{\pi}_1^{(k)}(1 - \widehat{\pi}_1^{(k)}), \dots, \widehat{\pi}_K^{(k)}(1 - \widehat{\pi}_K^{(k)}))$$

be an $N \times N$ matrix, where the i -th element is repeated m_i times, define

$$\widehat{Z}_{ij}^{(k)} = \mathbf{X}_i^\top \widehat{\boldsymbol{\beta}}^{(k)} + \frac{Y_{ij}^* - \widehat{\pi}_i^{(k)}}{\widehat{\pi}_i^{(k)}(1 - \widehat{\pi}_i^{(k)})},$$

and create an N -vector $\widehat{\mathbf{Z}}^{(k)} = (\widehat{Z}_{11}^{(k)}, \dots, \widehat{Z}_{Km_K}^{(k)})^\top$. The IWLS algorithm iterates

$$\widehat{\boldsymbol{\beta}}_n^{(k+1)} = (\mathbb{X}^\top \widehat{\mathbb{W}}^{(k)} \mathbb{X})^{-1} (\mathbb{X}^\top \widehat{\mathbb{W}}^{(k)} \widehat{\mathbf{Z}}^{(k)})$$

until convergence.

With the binomial format B, the regression matrix \mathbb{X}_R includes each observed value of \mathbf{X}_i only once, and its dimension is $K \times p$. Let

$$\hat{\mathbb{W}}_R^{(k)} = \text{diag} (m_1 \hat{\pi}_1^{(k)} (1 - \hat{\pi}_1^{(k)}), \dots, m_K \hat{\pi}_K^{(k)} (1 - \hat{\pi}_K^{(k)}))$$

be an $K \times K$ matrix, where each element appears just once, define

$$\hat{Z}_{Ri}^{(k)} = \mathbf{X}_i^\top \hat{\boldsymbol{\beta}}^{(k)} + \frac{Y_i - m_i \hat{\pi}_i^{(k)}}{m_i \hat{\pi}_i^{(k)} (1 - \hat{\pi}_i^{(k)})},$$

and create a K -vector $\hat{\mathbf{Z}}_R^{(k)} = (\hat{Z}_{R1}^{(k)}, \dots, \hat{Z}_{RK}^{(k)})^\top$. The IWLS algorithm iterates

$$\hat{\boldsymbol{\beta}}_n^{(k+1)} = (\mathbb{X}_R^\top \hat{\mathbb{W}}_R^{(k)} \mathbb{X}_R)^{-1} (\mathbb{X}_R^\top \hat{\mathbb{W}}_R^{(k)} \hat{\mathbf{Z}}_R^{(k)})$$

until convergence.

Obviously, $\mathbb{X}^\top \hat{\mathbb{W}}^{(k)} \mathbb{X} = \mathbb{X}_R^\top \hat{\mathbb{W}}_R^{(k)} \mathbb{X}_R$ and $\mathbb{X}^\top \hat{\mathbb{W}}^{(k)} \hat{\mathbf{Z}}^{(k)} = \mathbb{X}_R^\top \hat{\mathbb{W}}_R^{(k)} \hat{\mathbf{Z}}_R^{(k)}$, so the two implementations of the IWLS algorithm for the two data formats are equivalent.

The information matrix is

$$I(\boldsymbol{\beta}) = \mathbf{E}_{\mathbf{X}} \hat{\pi}_i (1 - \hat{\pi}_i) \mathbf{X}_i^{\otimes 2},$$

and it can be estimated by

$$\hat{I}_n = \frac{1}{N} \mathbb{X}^\top \hat{\mathbb{W}} \mathbb{X} = \frac{1}{N} \mathbb{X}_R^\top \hat{\mathbb{W}}_R \mathbb{X}_R.$$

The estimated variance of $\hat{\boldsymbol{\beta}}_n$ is

$$(\mathbb{X}^\top \hat{\mathbb{W}} \mathbb{X})^{-1} = (\mathbb{X}_R^\top \hat{\mathbb{W}}_R \mathbb{X}_R)^{-1}.$$

Deviance

The saturated model has K parameters π_1, \dots, π_K and their MLEs are $\tilde{\pi}_i = Y_i/m_i$. Under the first kind of asymptotics (K constant, $m_i \rightarrow \infty \forall i$), the MLE theory holds for the saturated model and the MLEs are consistent and asymptotically normal. For the second kind of asymptotics ($K \rightarrow \infty$, m_i small), the MLEs are inconsistent and the theory fails.

The deviance for logistic regression is

$$\begin{aligned} D(\mathbf{Y} \mid \hat{\boldsymbol{\beta}}) &= 2 \sum_{i=1}^K \sum_{j=1}^{m_i} \left[Y_{ij}^* \left(\log \frac{Y_i/m_i}{1 - Y_i/m_i} - \log \frac{\hat{\pi}_i}{1 - \hat{\pi}_i} \right) - \log \frac{1}{1 - Y_i/m_i} + \log \frac{1}{1 - \hat{\pi}_i} \right] \\ &= 2 \sum_{i=1}^K \left[Y_i \log \frac{Y_i}{m_i \hat{\pi}_i} + (m_i - Y_i) \log \frac{m_i - Y_i}{m_i - m_i \hat{\pi}_i} \right]. \end{aligned}$$

According to Theorem 2.9(iii), the difference in deviances between a submodel and a wider model has a limiting χ^2 distribution if the submodel holds. This is used for model building in logistic regression.

Under the first kind of asymptotics (K constant, $m_i \rightarrow \infty \forall i$), the deviance alone has asymptotically χ^2 distribution with $K - p$ degrees of freedom if the current model is valid. Thus, in these circumstances a formal goodness-of-fit test can be obtained by comparing the fit of the current model to the saturated model. If the deviance exceeds $1 - \alpha$ quantile of χ^2_{K-p} distribution, the current model should be rejected. However, this test is only valid if all covariates are discrete and there are enough observations at each value of \mathbf{X} (e.g., $m_i \hat{\pi}_i > 5$, $m_i(1 - \hat{\pi}_i) > 5$).

Pearson X^2

The Pearson residuals and the Pearson X^2 statistic have different versions and different interpretations depending on the data format. The two versions will be the same if and only if $m_i = 1$ for all i . It is important to keep in mind that most statistical software will calculate the version that corresponds to the input data format for the particular dataset.

Bernoulli data format

With the Bernoulli data format A, the Pearson residual is

$$r_{ij}^P = \frac{Y_{ij}^* - \hat{\pi}_i}{\sqrt{\hat{\pi}_i(1 - \hat{\pi}_i)}} = \begin{cases} \sqrt{\frac{1 - \hat{\pi}_i}{\hat{\pi}_i}} & Y_{ij}^* = 1 \\ -\sqrt{\frac{\hat{\pi}_i}{1 - \hat{\pi}_i}} & Y_{ij}^* = 0. \end{cases}$$

When plotted, these residuals form two separated clouds, one for successes and one for failures. Their usefulness for various model checking procedures is limited. The Pearson X^2 statistic is

$$X^2 = \sum_{i=1}^K \sum_{j=1}^{m_i} \frac{(Y_{ij}^* - \hat{\pi}_i)^2}{\hat{\pi}_i(1 - \hat{\pi}_i)} = \sum_{i=1}^K \left[\frac{Y_i}{m_i \hat{\pi}_i} m_i(1 - \hat{\pi}_i) + \frac{m_i - Y_i}{m_i(1 - \hat{\pi}_i)} m_i \hat{\pi}_i \right].$$

If the model fits well then $X^2 \approx N$. This version of the Pearson X^2 statistic does not possess a useful interpretation.

Binomial data format

With the binomial data format B, the Pearson residuals are calculated as $\hat{\mathbb{W}}_R^{1/2}(\hat{\mathbf{Z}}_R - \hat{\mathbb{X}}_R \hat{\boldsymbol{\beta}})$. The i -th residual is

$$r_i^P = \frac{Y_i - m_i \hat{\pi}_i}{\sqrt{m_i \hat{\pi}_i(1 - \hat{\pi}_i)}}$$

As long as $m_i \gg 1$, these residuals have much more useful distribution. They are

approximately normal as $m_i \rightarrow \infty$. The Pearson X^2 statistic is

$$X^2 = \sum_{i=1}^K \frac{(Y_i - m_i \hat{\pi}_i)^2}{m_i \hat{\pi}_i (1 - \hat{\pi}_i)} = \sum_{i=1}^K \frac{(Y_i - m_i \hat{\pi}_i)^2}{m_i \hat{\pi}_i} + \sum_{i=1}^K \frac{[m_i - Y_i - m_i(1 - \hat{\pi}_i)]^2}{m_i(1 - \hat{\pi}_i)},$$

which is the χ^2 test statistic for testing goodness of fit in a $2 \times K$ table with p estimated parameters. When K is constant and $m_i \rightarrow \infty$ for all i (at the same rate) then $X^2 \xrightarrow{D} \chi_{K-p}^2$ if the model is valid. For the saturated model ($p = K$), $X^2 = 0$. Thus, the binomial version of the Pearson X^2 statistic can be used for testing the goodness of fit if all covariates are discrete and there are enough observations at each value of \mathbf{X} . However, differences in Pearson X^2 statistics between the model and a submodel do not converge to a χ^2 distribution. This is why we always prefer deviance tests to tests based on the Pearson X^2 statistic.

3.2 Analysis of Poisson Count Data

3.2.1 Poisson loglinear model

Let Y_1, \dots, Y_n be independent random variables, $Y_i \sim \text{Po}(\lambda_i)$. Let λ_i depend on covariates \mathbf{X}_i through the identity $\log \lambda_i = \mathbf{X}_i^\top \boldsymbol{\beta}$.

Poisson distribution belongs to the exponential family with $\mu_i = \lambda_i$ and $\text{var } Y_i = \lambda_i$. The variance function is $V(\mu) = \mu$, the dispersion parameter is $\varphi = 1$, the canonical parameter is $\theta_i = \log \lambda_i$, $b(\theta_i) = e^{\theta_i}$. The log link is canonical for Poisson distribution.

We have

$$\mathbb{E}[Y_i | \mathbf{X}_i] = \text{var}[Y_i | \mathbf{X}_i] = \lambda_i = \exp\{\mathbf{X}_i^\top \boldsymbol{\beta}\}.$$

Interpretation of regression parameters

Let $\mathbf{X}^\top \boldsymbol{\beta} = \beta_1 + \beta_2 X_2 + \dots + \beta_p X_p$. Denote $\lambda_0 = \mathbb{E}[Y_i | X_2 = \dots = X_p = 0]$. Then $\log \lambda_0 = \beta_1$ so e^{β_1} is the expected value of Y_i for an individual with zero values in all covariates.

Now consider two individuals with observed covariates $\mathbf{x}^0 = (1, x_2, \dots, x_p)^\top$ and $\mathbf{x}^j = (1, x_2, \dots, x_j + 1, \dots, x_p)^\top$ (the j -th covariate is increased by 1, the others are the same). Denote $\lambda_{X0} = \mathbb{E}[Y_i | \mathbf{X} = \mathbf{x}^0]$ and $\lambda_{Xj} = \mathbb{E}[Y_i | \mathbf{X} = \mathbf{x}^j]$. Then

$$\beta_j = \log \frac{\lambda_{Xj}}{\lambda_{X0}} \quad \text{and} \quad e^{\beta_j} = \frac{\lambda_{Xj}}{\lambda_{X0}}.$$

Thus e^{β_j} is the proportional increase in $\mathbb{E} Y_i$ per unit difference in the covariate X_j . When $\beta_j = 0$ the ratio is 1 and the covariate has no effect on the expectation given the other covariates.

Estimation of parameters

The likelihood is

$$L_n(\boldsymbol{\beta} \mid \mathbf{Y}) = \prod_{i=1}^n \exp\{Y_i \log \lambda_i - \lambda_i - \log Y_i!\}$$

log-likelihood

$$\ell_n(\boldsymbol{\beta} \mid \mathbf{Y}) = \sum_{i=1}^n (Y_i \log \lambda_i - \lambda_i - \log Y_i!).$$

By Theorem 2.3, the score statistic with the canonical link is

$$\mathbf{U}_n(\boldsymbol{\beta} \mid \mathbf{Y}) = \sum_{i=1}^n (Y_i - \lambda_i) \mathbf{X}_i$$

and $\hat{\boldsymbol{\beta}}_n$ solves the equations

$$\sum_{i=1}^n Y_i \mathbf{X}_i = \sum_{i=1}^n \hat{\lambda}_i \mathbf{X}_i,$$

where

$$\hat{\lambda}_i = \exp\{\mathbf{X}_i^\top \hat{\boldsymbol{\beta}}_n\}.$$

The MLE of $\boldsymbol{\beta}$ is calculated by the IWLS algorithm

$$\hat{\boldsymbol{\beta}}_n^{(k+1)} = (\mathbb{X}^\top \hat{\mathbb{W}}^{(k)} \mathbb{X})^{-1} (\mathbb{X}^\top \hat{\mathbb{W}}^{(k)} \hat{\mathbf{Z}}^{(k)})$$

with

$$\hat{\mathbb{W}}^{(k)} = \text{diag}(\hat{\lambda}_1^{(k)}, \dots, \hat{\lambda}_n^{(k)})$$

and $\hat{\mathbf{Z}}^{(k)} = (\hat{Z}_1^{(k)}, \dots, \hat{Z}_n^{(k)})^\top$, where

$$\hat{Z}_i^{(k)} = \mathbf{X}_i^\top \hat{\boldsymbol{\beta}}^{(k)} + \frac{Y_i - \hat{\lambda}_i^{(k)}}{\hat{\lambda}_i^{(k)}}.$$

The information matrix is $I(\boldsymbol{\beta}) = \mathbf{E}_{\mathbf{X}} \lambda_i \mathbf{X}_i^{\otimes 2}$, which can be estimated by

$$\hat{I}_n = \frac{1}{n} \mathbb{X}^\top \hat{\mathbb{W}} \mathbb{X} = \frac{1}{n} \sum_{i=1}^n \hat{\lambda}_i \mathbf{X}_i^{\otimes 2}.$$

The estimated variance of $\hat{\boldsymbol{\beta}}_n$ is $(\mathbb{X}^\top \hat{\mathbb{W}} \mathbb{X})^{-1}$.

Deviance

The MLEs of the saturated model parameters are $\tilde{\lambda}_i = Y_i$. The deviance for loglinear model is

$$D(\mathbf{Y} | \hat{\boldsymbol{\beta}}) = 2 \sum_{i=1}^n \left[Y_i \log \frac{Y_i}{\hat{\lambda}_i} - (Y_i - \hat{\lambda}_i) \right]. \quad (3.2)$$

According to Theorem 2.9(iii), the difference in deviances between a submodel and a wider model has a limiting χ^2 distribution if the submodel holds. This is used for loglinear model building.

Pearson X^2

The Pearson residual is

$$r_i^P = \frac{Y_i - \hat{\lambda}_i}{\sqrt{\hat{\lambda}_i}},$$

the Pearson X^2 statistic is

$$X^2 = \sum_{i=1}^n \frac{(Y_i - \hat{\lambda}_i)^2}{\hat{\lambda}_i}. \quad (3.3)$$

3.2.2 Modelling Poisson process intensity

Let Y_i be a number of events observed during t_i units of time on a certain individual (experimental unit), let \mathbf{X}_i be covariates. We have independent realizations of (Y_i, t_i, \mathbf{X}_i) for $i = 1, \dots, n$.

Suppose $Y_i \sim \text{Po}(\lambda_i t_i)$, that is, the outcomes are generated by a homogeneous Poisson process with intensity λ_i . Let λ_i depend on the covariates \mathbf{X}_i through the identity $\log \lambda_i = \mathbf{X}_i^T \boldsymbol{\beta}$. Then

$$\mathbf{E}[Y_i | \mathbf{X}_i] = \text{var}[Y_i | \mathbf{X}_i] = t_i \lambda_i = t_i \exp\{\mathbf{X}_i^T \boldsymbol{\beta}\} = \exp\{\log t_i + \beta_1 + \beta_2 X_{i2} + \dots + \beta_p X_{ip}\}.$$

The intensity $\lambda_i = \mathbf{E} Y_i / t_i$ describes the expected number of events observed during a unit time interval.

This model can be fitted as a Poisson loglinear model with the term $\log t_i$ added to the linear predictor. It is a term similar to a covariate, except that its parameter is equal to 1 by default and is not estimated. Such a term is called *an offset** in the GLM terminology. Estimation and testing proceeds in exactly the same way as if offset was not present.

* Česky *offset*

3.3 Loglinear Models for Contingency Tables

3.3.1 Two-way contingency table

Consider discrete random variables $X \in \{1, \dots, I\}$ and $Z \in \{1, \dots, J\}$. Observe n independent realizations $(X_1, Z_1), \dots, (X_n, Z_n)$ of this pair. Denote the observed count of the pair $(X = i, Z = j)$ by $n_{ij} = \sum_{l=1}^n I(X_l = i, Z_l = j)$. The observed counts (also called *frequencies*) can be arranged into a two-way contingency table

	$Z = 1$	\dots	$Z = J$	Total
$X = 1$	n_{11}	\dots	n_{1J}	n_{1+}
\vdots	\vdots		\vdots	\vdots
$X = I$	n_{I1}	\dots	n_{IJ}	n_{I+}
Total	n_{+1}	\dots	n_{+J}	$n_{++} = n$

where $n_{i+} = \sum_{j=1}^J n_{ij}$ and $n_{+j} = \sum_{i=1}^I n_{ij}$.

Denote the expected cell frequencies

$$m_{ij} = \mathbb{E} n_{ij}, \quad m_{i+} = \sum_{j=1}^J m_{ij}, \quad m_{+j} = \sum_{i=1}^I m_{ij},$$

and $m_{++} = \sum_{i=1}^I \sum_{j=1}^J m_{ij}$. The cell probabilities are be

$$\pi_{ij} = \mathbb{P}[X = i, Z = j], \quad \pi_{i+} = \sum_{j=1}^J \pi_{ij} = \mathbb{P}[X = i], \quad \pi_{+j} = \sum_{i=1}^I \pi_{ij} = \mathbb{P}[Z = j].$$

Obviously, $\pi_{++} = \sum_{i=1}^I \sum_{j=1}^J \pi_{ij} = 1$.

The expected frequencies are related to the cell probabilities as follows (we allow n to be random):

$$m_{ij} = \mathbb{E} n_{ij} = \mathbb{E} \sum_{l=1}^n I(X_l = i, Z_l = j) = \mathbb{E} \mathbb{E} \left[\sum_{l=1}^n I(X_l = i, Z_l = j) \mid n \right] = \mathbb{E} n \pi_{ij} = m_{++} \pi_{ij},$$

so $\pi_{ij} = m_{ij}/m_{++}$.

The goal is to use the observed counts n_{ij} to model the cell probabilities π_{ij} , investigate the marginal distributions of X and Z and the associations between X and Z .

3.3.2 Distributions of observed counts

There are several distributions we can use for observed counts in the contingency table.

Poisson distribution

Let n_{11}, \dots, n_{IJ} be independent random variables with Poisson distributions $n_{ij} \sim \text{Po}(m_{ij})$. It follows $\mathbb{E} n_{ij} = \text{var} n_{ij} = m_{ij}$. The joint density of the whole table is

$$P [n_{11} = k_{11}, \dots, n_{IJ} = k_{IJ}] = \prod_{i,j} \frac{1}{k_{ij}!} m_{ij}^{k_{ij}} e^{-m_{ij}}, \quad k_{ij} = 0, 1, 2, \dots$$

The total number of observations $n = \sum_{i,j} n_{ij}$ is a random variable with the distribution $\text{Po}(m_{++})$.

The asymptotics does not work by observing an increasing number of independent Poisson variables (the total IJ is fixed) but by letting $m_{++} \rightarrow \infty$. The asymptotic MLE theory for iid data does not apply to this case.

Multinomial distribution

Let the vector (n_{11}, \dots, n_{IJ}) follow the multinomial distribution $\text{Mult}_{IJ}(n, \boldsymbol{\pi})$, where $\boldsymbol{\pi} = (\pi_{11}, \dots, \pi_{IJ})^\top$. The joint density of the whole table is

$$P [n_{11} = k_{11}, \dots, n_{IJ} = k_{IJ}] = n! \prod_{i,j} \frac{1}{k_{ij}!} \pi_{ij}^{k_{ij}}, \quad k_{ij} = 0, 1, \dots, n, \quad \sum_{i,j} k_{ij} = n.$$

The total number of observations n is fixed, $n_{ij} \sim \text{Bi}(n, \pi_{ij})$, $\mathbb{E} n_{ij} = n\pi_{ij}$, $\text{var} n_{ij} = n\pi_{ij}(1 - \pi_{ij})$, the counts are not independent.

The contingency table can be expressed by summing n iid random vectors, each with distribution $\text{Mult}_{IJ}(1, \boldsymbol{\pi})$. The asymptotics works through letting $n \rightarrow \infty$. The asymptotic MLE theory for iid data applies to this case.

Row multinomial distribution

Let the vectors (n_{i1}, \dots, n_{iJ}) , $i = 1, \dots, I$, be iid with the multinomial distribution $\text{Mult}_J(n_{i+}, \boldsymbol{\pi}_i)$, where $\boldsymbol{\pi}_i = (\pi_{i1}/\pi_{i+}, \dots, \pi_{iJ}/\pi_{i+})^\top$. The joint density of the whole table is

$$P [n_{11} = k_{11}, \dots, n_{IJ} = k_{IJ}] = \prod_i n_{i+}! \prod_j \frac{1}{k_{ij}!} \left(\frac{\pi_{ij}}{\pi_{i+}} \right)^{k_{ij}}, \quad k_{ij} = 0, 1, \dots, n, \quad \sum_j k_{ij} = n_{i+}.$$

The numbers of observations n_{i+} in the I rows of the table are fixed, $n_{ij} \sim \text{Bi}(n_{i+}, \frac{\pi_{ij}}{\pi_{i+}})$, $\mathbb{E} n_{ij} = n_{i+} \frac{\pi_{ij}}{\pi_{i+}}$, $\text{var} n_{ij} = n_{i+} \frac{\pi_{ij}}{\pi_{i+}} (1 - \frac{\pi_{ij}}{\pi_{i+}})$, the counts are independent between rows but dependent within rows.

The asymptotics works through letting $n_{i+} \rightarrow \infty$ for all i at the same rate. The asymptotic MLE theory for iid data applies to this case.

Equivalence of Poisson and multinomial models

We start with a result stating that Poisson and multinomial distributions are related through conditioning on the total count.

Lemma 3.1. *Let $X_i \sim \text{Po}(\lambda_i)$ be independent random variables, $i = 1, \dots, n$. Then the conditional joint distribution of the random vector $(X_1, \dots, X_n)^\top$ given $\sum_{i=1}^n X_i = s$ is $\text{Mult}_n(s, \mathbf{p})$, where $\mathbf{p} = (p_1, \dots, p_n)^\top$ and $p_i = \lambda_i / \sum_{j=1}^n \lambda_j$.*

Corollary. Let $n_{ij} \sim \text{Po}(m_{ij})$ be independent, $i = 1, \dots, I$, $j = 1, \dots, J$. Then:

- The conditional joint distribution of $(n_{11}, \dots, n_{IJ})^\top$ given $n_{++} = n$ is $\text{Mult}_{IJ}(n, \boldsymbol{\pi})$, where the components of $\boldsymbol{\pi}$ are $\pi_{ij} = m_{ij}/m_{++}$ (\Rightarrow multinomial distribution).
- The conditional joint distribution of $(n_{i1}, \dots, n_{iJ})^\top$ given n_{i+} is $\text{Mult}_J(n_{i+}, \boldsymbol{\pi}_i)$, where the components of $\boldsymbol{\pi}_i$ are $\pi_{ij} = m_{ij}/m_{i+} = \pi_{ij}/\pi_{i+}$ (\Rightarrow row multinomial distribution).

The corollary states that both multinomial and row multinomial distributions of a contingency table can be obtained from Poisson distribution by conditioning on some observed totals.

Assume that the loglinear model holds for the expected frequencies m_{ij} , in particular,

$$\log \mathbb{E} n_{ij} = \alpha + \boldsymbol{\beta}^\top \mathbf{X}_{ij} \quad \text{or} \quad m_{ij} = e^{\alpha + \boldsymbol{\beta}^\top \mathbf{X}_{ij}},$$

where α is the intercept and \mathbf{X}_{ij} is a vector of covariates characterizing the (i, j) -th cell. The maximum dimension of \mathbf{X}_{ij} is $IJ - 1$. The cell probabilities π_{ij} can be expressed as follows:

$$\pi_{ij} = \frac{m_{ij}}{m_{++}} = \frac{e^{\boldsymbol{\beta}^\top \mathbf{X}_{ij}}}{\sum_{k,l} e^{\boldsymbol{\beta}^\top \mathbf{X}_{kl}}}. \quad (3.4)$$

Theorem 3.2. *The likelihood functions for estimation of parameters $\boldsymbol{\beta}$ in the loglinear model $\log m_{ij} = \alpha + \boldsymbol{\beta}^\top \mathbf{X}_{ij}$ arising from Poisson or multinomial sampling distributions are equivalent (they differ only by a multiplicative constant that does not depend on $\boldsymbol{\beta}$).*

Note. Theorem 3.2 does not deal with estimation of the intercept α – in fact, the intercept is not even identifiable in the multinomial model. This is obvious from expression (3.4).

Theorem 3.2 can be extended to row multinomial distribution as follows:

Theorem 3.3. (*Palmgren 1981*) *The likelihood functions for estimation of parameters $\boldsymbol{\beta}$ in the loglinear model $\log m_{ij} = \alpha_i + \boldsymbol{\beta}^\top \mathbf{X}_{ij}$ arising from Poisson or row-multinomial sampling distributions are equivalent (they differ only by a multiplicative constant that does not depend on $\boldsymbol{\beta}$).*

Note. Row multinomial sampling requires row-specific intercept in the loglinear model.

Corollary. Expressions for any quantity derived from the likelihood function for β (score function, information matrix, the MLE) and their properties (asymptotic distributions, test statistics, confidence intervals) are the same no matter which of the three distributions generated the contingency table.

When the data are generated by the Poisson model, they can be transformed to the multinomial model by conditioning on the observed cell count total $n = n_{++}$. The asymptotic results hold in the multinomial model (the data are equivalent to n independent observations from Mult_1). The formulae can be derived from the theory of GLM for the loglinear model with Poisson distribution.

In the rest of this section we assume the loglinear model with Poisson distribution but the results apply to multinomial models without change.

3.3.3 Loglinear models for two-way tables

The independence model (X, Z)

The model (X, Z) for expected cell counts is defined by the equation

$$\log m_{ij} = \alpha + \beta_i^X + \beta_j^Z \quad (3.5)$$

with the constraints $\beta_1^X = \beta_1^Z = 0$. The regressors that generate this model are

$$\begin{aligned} \mathbf{X}_{11} &= (1, 0, \dots, 0, \dots, 0, \dots, 0, \dots, 0)^\top && \text{for the cell } (1, 1), \\ \mathbf{X}_{i1} &= (1, 0, \dots, 1, \dots, 0, \dots, 0, \dots, 0)^\top && \text{for the cell } (i, 1), \quad i \neq 1, \\ \mathbf{X}_{1j} &= (1, 0, \dots, 0, \dots, 0, \dots, 1, \dots, 0)^\top && \text{for the cell } (1, j), \quad j \neq 1, \\ \mathbf{X}_{ij} &= (1, 0, \dots, 1, \dots, 0, \dots, 1, \dots, 0)^\top && \text{for the cell } (i, j), \quad i \neq 1, \quad j \neq 1. \end{aligned}$$

The parameters are

$$\beta = (\alpha, \beta_2^X, \dots, \beta_I^X, \beta_2^Z, \dots, \beta_J^Z)^\top.$$

the dimension of the parameter vector is $p = 1 + I - 1 + J - 1 = I + J - 1$.

The parameter α can be expressed as

$$\alpha = \log m_{11} = \log m_{++} + \log \pi_{11},$$

hence the model can be also stated in terms of cell probabilities

$$\log \pi_{ij} = \log \frac{m_{ij}}{m_{++}} = \log \pi_{11} + \beta_i^X + \beta_j^Z. \quad (3.6)$$

The interpretation of the regression parameters is as follows:

$$e^{\beta_i^X} = \frac{\pi_{i+}}{\pi_{1+}} = \frac{\text{P}[X = i]}{\text{P}[X = 1]} \quad \text{and} \quad e^{\beta_j^Z} = \frac{\pi_{+j}}{\pi_{+1}} = \frac{\text{P}[Z = j]}{\text{P}[Z = 1]},$$

that is, $e^{\beta_i^X}$ is the odds of observing the i -th level of the variable X compared to observing the first level.

The cell probabilities satisfy the equation

$$\pi_{ij} = \pi_{i+}\pi_{+j}$$

for all i and j , which means that the model (X, Z) holds if and only if the variables X and Z are independent.

The fitted cell frequencies in this model can be expressed explicitly

$$\hat{m}_{ij} = n\hat{\pi}_{ij} = n\hat{\pi}_{i+}\hat{\pi}_{+j} = n\frac{n_{i+}}{n}\frac{n_{+j}}{n} = \frac{n_{i+}n_{+j}}{n}$$

because $\frac{n_{i+}}{n}$ and $\frac{n_{+j}}{n}$ are the MLE's of π_{i+} and π_{+j} .

The interaction model (XZ)

The model (XZ) for expected cell counts is defined by the equation

$$\log m_{ij} = \alpha + \beta_i^X + \beta_j^Z + \beta_{ij}^{XZ} \quad (3.7)$$

with the constraints $\beta_1^X = \beta_1^Z = 0$, $\beta_{i1}^{XZ} = 0$ for all $i = 1, \dots, I$, and $\beta_{1j}^{XZ} = 0$ for all $j = 1, \dots, J$. The number of parameters in this model is $p = I + J - 1 + (I - 1)(J - 1) = IJ$. This is the saturated model, hence the estimated expected cell counts \hat{m}_{ij} (fitted values) are equal to the observed cell counts n_{ij} for all i, j .

The equivalent model for cell probabilities is

$$\log \pi_{ij} = \log \pi_{11} + \beta_i^X + \beta_j^Z + \beta_{ij}^{XZ}. \quad (3.8)$$

The interpretation of the main effects is as follows:

$$e^{\beta_i^X} = \frac{\pi_{i1}}{\pi_{11}} = \frac{P(X = i | Z = 1)}{P(X = 1 | Z = 1)} \quad \text{and} \quad e^{\beta_j^Z} = \frac{\pi_{1j}}{\pi_{11}} = \frac{P(X = 1 | Z = j)}{P(X = 1 | Z = 1)},$$

that is, $e^{\beta_i^X}$ is the odds of observing the i -th level of the variable X compared to observing the first level of X when $Z = 1$.

The interaction parameters can be written as

$$e^{\beta_{ij}^{XZ}} = \frac{\pi_{ij}\pi_{11}}{\pi_{i1}\pi_{1j}} = \frac{P[X = i, Z = j]P[X = 1, Z = 1]}{P[X = i, Z = 1]P[X = 1, Z = j]} = \frac{P(X = i | Z = j)P(X = 1 | Z = 1)}{P(X = 1 | Z = j)P(X = i | Z = 1)},$$

which is the odds ratio in the 2×2 sub-table that includes the first and the i -th rows and the first and j -th columns from the original table. The odds ratio expresses the

proportional change in the odds of the event $X = i$ (relative to $X = 1$) when Z changes from 1 to j .

The cell probabilities can be expressed as

$$\pi_{ij} = P[X = i, Z = j] = \frac{e^{\beta_i^X + \beta_j^Z + \beta_{ij}^{XZ}}}{\sum_{k,l} e^{\beta_k^X + \beta_l^Z + \beta_{kl}^{XZ}}}.$$

Since this is the saturated model, its deviance $D(XZ)$ is zero. The MLE theory holds for this model because all covariates are discrete and the number of parameters $p = IJ$ is constant.

3.3.4 Testing independence in a two-way table

The variables X and Z are independent if and only if the model (X, Z) holds, that is, all the interaction parameters are zero. So, a test of independence should test all $(I - 1)(J - 1)$ interaction parameters simultaneously.

Such a test can be based on the deviance $D(X, Z)$ of the independence model. From (3.2), the deviance is

$$D(X, Z) = 2 \sum_{i=1}^I \sum_{j=1}^J \left[n_{ij} \log \frac{n_{ij}}{\hat{m}_{ij}} - (n_{ij} - \hat{m}_{ij}) \right],$$

where $\hat{m}_{ij} = n_{i+}n_{+j}/n$. If independence holds then $D(X, Z) \xrightarrow{D} \chi_{(I-1)(J-1)}^2$, so the hypothesis is rejected at the level of α when $D(X, Z) \geq \chi_{(I-1)(J-1)}^2(1 - \alpha)$.

The independence hypothesis can be also tested by the Pearson X^2 statistic. By (3.3), the Pearson X^2 statistic is

$$X^2 = \sum_{i=1}^I \sum_{j=1}^J \frac{(n_{ij} - \hat{m}_{ij})^2}{\hat{m}_{ij}}.$$

This is the classical χ^2 statistic for testing independence in a two-way table. Under the null hypothesis, it also converges in distribution to $\chi_{(I-1)(J-1)}^2$.

3.3.5 Loglinear models for three-way tables

Three-way contingency table

Consider categorical random variables $X \in \{1, \dots, I\}$, $Z \in \{1, \dots, J\}$, and $V \in \{1, \dots, K\}$. Observe n independent realizations $(X_1, Z_1, V_1), \dots, (X_n, Z_n, V_n)$ of this triplet. Denote

the observed count of the outcome $(X = i, Z = j, V = k)$ by

$$n_{ijk} = \sum_{l=1}^n I(X_l = i, Z_l = j, V_l = k).$$

The observed counts n_{ijk} form a three-way contingency table.

Let $n_{ij+} = \sum_{k=1}^K n_{ijk}$, $n_{i++} = \sum_{j=1}^J \sum_{k=1}^K n_{ijk}$, etc. Denote $m_{ijk} = \mathbb{E} n_{ijk}$ and $\pi_{ijk} = \mathbb{P}[X = i, Z = j, V = k]$. The symbols m_{ij+} , m_{++k} , π_{i+k} , π_{+j+} etc. all have the obvious meaning (summing over the indices replaced by +). Obviously, $\pi_{+++} = \sum_{i,j,k} \pi_{ijk} = 1$.

Marginal and conditional associations in a three-way table

Suppose we would like to describe the associations between the variables X and Z in the presence of a third classifier V . We have two ways to do this:

Marginal associations

Let us ignore the existence of the third variable. We can form a two-way contingency table for X and Z by collapsing the original three-way table n_{ijk} across the levels of V , getting a two-way table n_{ij+} . Now the associations between X and Z can be described and investigated by the methods for a two-way table described in Section 3.3.3 applied to n_{ij+} .

Definition 3.1. Discrete variables X and Z are called *marginally independent* (in the presence of a third discrete variable V) if and only if

$$\pi_{ij+} = \pi_{i++}\pi_{+j+} \quad \text{for all } i, j.$$

Marginal associations between X and Z are described by marginal odds ratios, which correspond to exponentiated interaction terms in the model (3.8) applied to the collapsed two-way table.

Definition 3.2. The marginal odds ratio for the i -th level of X and the j -th level of Z are defined as

$$\theta_{ij}^{XZ} = \frac{\pi_{ij+}\pi_{11+}}{\pi_{i1+}\pi_{1j+}}.$$

X and Z are marginally independent if and only if $\theta_{ij}^{XZ} = 1$ for all i and j .

Conditional associations

Let us fix the value of variable V at some value $k \in \{1, \dots, K\}$. We get K separate two-way contingency tables for X and Z by considering each of the layers of the three-way table n_{ijk} formed by fixing $V = k$. The k -th two-way table has observed counts

n_{ijk} ($i = 1, \dots, I, j = 1, \dots, J$) and cell probabilities proportional to π_{ijk} . The associations between X and Z can be described and investigated by the methods described in Section 3.3.3 applied to n_{ijk} for each fixed k .

Definition 3.3. Variables X and Z are called *conditionally independent given V* if and only if

$$P(X = i, Z = j | V = k) = P(X = i | V = k)P(Z = j | V = k) \quad \text{for all } i, j, k.$$

Theorem 3.4. Variables X and Z are conditionally independent given V if and only if

$$\pi_{ijk} = \frac{\pi_{i+k}\pi_{+jk}}{\pi_{++k}} \quad \text{for all } i, j, k.$$

Conditional associations between X and Z are described by conditional odds ratios.

Definition 3.4. The conditional odds ratios for the i -th level of X and the j -th level of Z given the k -th level of V are defined as

$$\theta_{ij(k)}^{XZ} = \frac{\pi_{ijk}\pi_{11k}}{\pi_{i1k}\pi_{1jk}}.$$

X and Z are conditionally independent if and only if $\theta_{ij(k)}^{XZ} = 1$ for all i, j , and k .

Comments

Marginal and conditional associations can be quite different. In each application, the analyst must decide whether marginal or conditional associations are more relevant for the particular case. The decision is based on the context of the problem and the desired interpretation of the results.

The variable V usually acts as a *confounder*. A confounder is a variable that affects both X and Z and thus masks the true relationship between X and Z . In the conditional analysis, we adjust for the confounding effect of V and reveal the true association of X and Z . Therefore, the conditional associations are the correct approach in this case.

However, it may happen that V is a part of the causal pathway between X and Z (X affects V and V affects Z). Then V cannot be considered a confounder. The conditional analysis of the association between X and Z removes the part of the effect that is mediated through V and thus gives misleading results about the true association. In this case, the marginal associations provide the correct approach to the problem.

The independence model (X, Z, V)

The model (X, Z, V) is defined by the equation

$$\log m_{ijk} = \alpha + \beta_i^X + \beta_j^Z + \beta_k^V \tag{3.9}$$

with the constraints $\beta_1^X = \beta_1^Z = \beta_1^V = 0$.

In terms of cell probabilities,

$$\log \pi_{ijk} = \log \pi_{111} + \beta_i^X + \beta_j^Z + \beta_k^V. \quad (3.10)$$

The dimension of the parameter vector is $p = I + J + K - 2$.

The interpretation of the regression parameters is as follows:

$$e^{\beta_i^X} = \frac{\pi_{i++}}{\pi_{1++}} = \frac{P[X = i]}{P[X = 1]}, \quad e^{\beta_j^Z} = \frac{\pi_{+j+}}{\pi_{+1+}} = \frac{P[Z = j]}{P[Z = 1]}, \quad e^{\beta_k^V} = \frac{\pi_{++k}}{\pi_{++1}} = \frac{P[V = k]}{P[V = 1]}.$$

The cell probabilities satisfy the equation

$$\pi_{ijk} = \pi_{i++}\pi_{+j+}\pi_{++k}$$

for all i, j, k . The model (X, Z, V) holds if and only if the variables X , Z and V are both marginally and conditionally independent.

Model (XV, Z)

The model (XV, Z) for expected cell counts is defined by the equation

$$\log m_{ijk} = \alpha + \beta_i^X + \beta_j^Z + \beta_k^V + \beta_{ik}^{XV} \quad (3.11)$$

with the constraints $\beta_1^X = \beta_1^Z = \beta_1^V = 0$, $\beta_{i1}^{XV} = 0$ for all i , and $\beta_{1k}^{XV} = 0$ for all k . The number of parameters in this model is $p = IK + J - 1$.

The model for cell probabilities is

$$\log \pi_{ijk} = \log \pi_{111} + \beta_i^X + \beta_j^Z + \beta_k^V + \beta_{ik}^{XV}. \quad (3.12)$$

The interpretation of the main effects:

$$e^{\beta_i^X} = \frac{\pi_{i+1}}{\pi_{1+1}} = \frac{P(X = i | V = 1)}{P(X = 1 | V = 1)}, \quad e^{\beta_k^V} = \frac{\pi_{1+k}}{\pi_{1+1}} = \frac{P(V = k | X = 1)}{P(V = 1 | X = 1)},$$

$$e^{\beta_j^Z} = \frac{\pi_{+j+}}{\pi_{+1+}} = \frac{P[Z = j]}{P[Z = 1]}.$$

It can be easily shown that

$$e^{\beta_{ik}^{XV}} = \theta_{ik}^{XV} = \theta_{ik(j)}^{XV},$$

which is both the marginal and the conditional odds ratio between X and V .

The cell probabilities satisfy the equation

$$\pi_{ijk} = \pi_{i+k}\pi_{+j+}$$

for all i, j, k . Hence in this model,

- the pair (X, V) is independent of Z ,
- X is both marginally and conditionally independent of Z ,
- X is associated with V , and
- the marginal and conditional associations between X and V agree.

Model (XV, ZV)

The model (XV, ZV) for expected cell counts is defined by the equation

$$\log m_{ijk} = \alpha + \beta_i^X + \beta_j^Z + \beta_k^V + \beta_{ik}^{XV} + \beta_{jk}^{ZV} \quad (3.13)$$

with the constraints $\beta_1^X = \beta_1^Z = \beta_1^V = 0$, $\beta_{i1}^{XV} = 0$ for all i , $\beta_{1k}^{XV} = 0$ for all k , $\beta_{j1}^{ZV} = 0$ for all j , and $\beta_{1k}^{ZV} = 0$ for all k . The number of parameters in this model is $p = K(I + J - 1)$.

The model for cell probabilities is

$$\log \pi_{ijk} = \log \pi_{111} + \beta_i^X + \beta_j^Z + \beta_k^V + \beta_{ik}^{XV} + \beta_{jk}^{ZV}. \quad (3.14)$$

The interpretation of the main effects:

$$e^{\beta_i^X} = \frac{\pi_{i+1}}{\pi_{1+1}} = \frac{P(X = i | V = 1)}{P(X = 1 | V = 1)}, \quad e^{\beta_j^Z} = \frac{\pi_{+j1}}{\pi_{+11}} = \frac{P(Z = j | V = 1)}{P(Z = 1 | V = 1)},$$

$$e^{\beta_k^V} = \frac{\pi_{11k}}{\pi_{111}} = \frac{P(V = k | X = 1, Z = 1)}{P(V = 1 | X = 1, Z = 1)}.$$

The cell probabilities satisfy the equation

$$\pi_{ijk} = \frac{\pi_{i+k}\pi_{+jk}}{\pi_{++k}}$$

for all i, j, k .

The association between X and Z :

$$\begin{aligned} \text{Conditional OR for } X \leftrightarrow Z: & \quad \theta_{ij(k)}^{XZ} = \frac{\pi_{ijk}\pi_{11k}}{\pi_{i1k}\pi_{1jk}} = 1 \\ \text{Marginal OR for } X \leftrightarrow Z: & \quad \theta_{ij}^{XZ} = \frac{\pi_{ij+}\pi_{11+}}{\pi_{i1+}\pi_{1j+}} \neq 1 \end{aligned}$$

The association between X and V :

$$\begin{aligned} \text{Conditional OR for } X \leftrightarrow V: & \quad \theta_{ik(j)}^{XV} = \frac{\pi_{ijk}\pi_{1j1}}{\pi_{ij1}\pi_{1jk}} = e^{\beta_{ik}^{XV}} \\ \text{Marginal OR for } X \leftrightarrow V: & \quad \theta_{ik}^{XV} = \frac{\pi_{i+k}\pi_{1+1}}{\pi_{i+1}\pi_{1+k}} = e^{\beta_{ik}^{XV}} \end{aligned}$$

The association between Z and V is similar to that between X and V (rotate indices). In particular,

$$\theta_{jk(i)}^{ZV} = \theta_{jk}^{ZV} = e^{\beta_{jk}^{ZV}}.$$

Hence in this model,

- all three variables are marginally dependent,
- X is conditionally independent of Z given V (but marginally dependent),
- X is associated with V and the marginal and conditional associations between X and V agree,
- Z is associated with V and the marginal and conditional associations between Z and V agree.

Model (XV, ZV, XZ)

The model (XV, ZV, XZ) for expected cell counts is defined by the equation

$$\log m_{ijk} = \alpha + \beta_i^X + \beta_j^Z + \beta_k^V + \beta_{ik}^{XV} + \beta_{jk}^{ZV} + \beta_{ij}^{XZ} \quad (3.15)$$

with the constraints $\beta_1^X = \beta_1^Z = \beta_1^V = 0$, $\beta_{i1}^{XV} = 0$ for all i , $\beta_{1k}^{XV} = 0$ for all k , $\beta_{j1}^{ZV} = 0$ for all j , $\beta_{1k}^{ZV} = 0$ for all k , $\beta_{i1}^{XZ} = 0$ for all i , and $\beta_{1j}^{XZ} = 0$ for all j . The number of parameters in this model is $p = KI + KJ + IJ - (I + J + K) + 1$.

The model for cell probabilities is

$$\log \pi_{ijk} = \log \pi_{111} + \beta_i^X + \beta_j^Z + \beta_k^V + \beta_{ik}^{XV} + \beta_{jk}^{ZV} + \beta_{ij}^{XZ} \quad (3.16)$$

The interpretation of the main effects:

$$e^{\beta_i^X} = \frac{\pi_{i11}}{\pi_{111}}, \quad e^{\beta_j^Z} = \frac{\pi_{1j1}}{\pi_{111}}, \quad e^{\beta_k^V} = \frac{\pi_{11k}}{\pi_{111}}.$$

The interaction parameters determine the conditional odds ratios:

$$\begin{aligned} \text{Conditional OR for } X \leftrightarrow Z: & \quad \theta_{ij(k)}^{XZ} = e^{\beta_{ij}^{XZ}} \\ \text{Conditional OR for } X \leftrightarrow V: & \quad \theta_{ik(j)}^{XV} = e^{\beta_{ik}^{XV}} \\ \text{Conditional OR for } Z \leftrightarrow V: & \quad \theta_{jk(i)}^{ZV} = e^{\beta_{jk}^{ZV}} \end{aligned}$$

Marginal odds ratios are all different from conditional odds ratios and do not have nice expressions. There is no decomposition of the cell probabilities π_{ijk} .

In this model,

- all three variables are marginally and conditionally dependent,
- marginal associations are different from conditional associations,
- conditional associations do not depend on the value of the conditioning variable.

Model (XZV)

The model (XZV) for expected cell counts is defined by the equation

$$\log m_{ijk} = \alpha + \beta_i^X + \beta_j^Z + \beta_k^V + \beta_{ik}^{XV} + \beta_{jk}^{ZV} + \beta_{ij}^{XZ} + \beta_{ijk}^{XZV} \quad (3.17)$$

with the usual constraints (any parameter with 1 anywhere among the indices is set to 0). The number of parameters in this model is $p = IJK$. It is the saturated model.

The model for cell probabilities is

$$\log \pi_{ijk} = \log \pi_{111} + \beta_i^X + \beta_j^Z + \beta_k^V + \beta_{ik}^{XV} + \beta_{jk}^{ZV} + \beta_{ij}^{XZ} + \beta_{ijk}^{XZV} \quad (3.18)$$

The interpretation of the main effects is the same as in the model (XV, ZV, XZ). The second-order interaction parameters determine the conditional odds ratios at the first level of the conditioning variable:

Conditional OR for $X \leftrightarrow Z$ given $V = 1$:	$\theta_{ij(1)}^{XZ} = e^{\beta_{ij}^{XZ}}$
Conditional OR for $X \leftrightarrow V$ given $Z = 1$:	$\theta_{ik(1)}^{XV} = e^{\beta_{ik}^{XV}}$
Conditional OR for $Z \leftrightarrow V$ given $X = 1$:	$\theta_{jk(1)}^{ZV} = e^{\beta_{jk}^{ZV}}$.

The third-order interaction parameters determine how the conditional odds ratios change if the conditioning is not on the first level:

$$e^{\beta_{ijk}^{XZV}} = \frac{\theta_{ij(k)}^{XZ}}{\theta_{ij(1)}^{XZ}} = \frac{\theta_{ik(j)}^{XV}}{\theta_{ik(1)}^{XV}} = \frac{\theta_{jk(i)}^{ZV}}{\theta_{jk(1)}^{ZV}}.$$

Marginal odds ratios are all different from conditional odds ratios and do not have nice expressions. There is no decomposition of the cell probabilities π_{ijk} .

This model is similar to (XV, ZV, XZ) except that conditional associations do depend on the value of the conditioning variable.

Model selection

The model selection strategy that usually works best is to fit the saturated model (XZV) first and then test whether the higher order interactions can be removed, starting with the three-way interaction. Model selection should be hierarchical (do not remove lower order terms if a higher order term involving the same variable is still in the model). Deviance tests are used to decide whether a term can be removed. Because the saturated model satisfies the MLE theory assumptions, the deviance of each fitted model can be compared with a quantile of a suitable χ^2 distribution to check whether the current model fits well.

The final model that cannot be further reduced reveals the structure of the associations between the three variables, as explained above.

3.3.6 Loglinear models for multi-way tables

Theoretically, a loglinear model can fit contingency tables of arbitrarily high dimension. The interpretation of the main effects, two-way and three-way interaction remains the same as in a three-way table, unless there are interactions of even higher order.

The dependence structure among multiple categorical variables can be deduced from an undirected graph where each variable plays the role of a node and two-way interactions are the edges between the nodes. If there is at least one path connecting two nodes, then the two variables corresponding to the nodes are marginally dependent. If all the paths connecting the two nodes can be interrupted by removing a certain set of nodes from the graph, then the two variables are conditionally independent given the set of variables corresponding to the removed nodes.

Notes on model selection for multi-way tables

- Model building proceeds by performing deviance tests comparing a model vs. a submodel.
- The number of possible models in a multi-way table can be very large. There is no way to fit all of them.
- The starting model cannot be too complex. In most cases, we cannot start by the saturated model. Interactions of the fourth or higher orders are difficult to interpret and we try to avoid them even if they seem to be statistically significant. A reasonable strategy could be to start with a model containing all three-way interactions, test its goodness-of-fit using deviance, and remove the insignificant three-way (and then two-way) interactions in a backward step-wise procedure. It is better to do this interactively rather than to use some automated model-building procedure.
- Multi-way tables require a lot of data. If there are four or five factors with a moderate number of levels, the number of cells in the table is huge and the observed counts can be quite low even if the total number of observations is in the thousands. If too many of the fitted cell counts are below 5, the asymptotic approximations tend to be unreliable. The analyst must make sure that there are enough observations in the cells, otherwise some of the variables must be removed from the analysis or their levels must be merged to reduce the number of cells in the table.

3.3.7 Equivalence of loglinear and logistic models

Consider three categorical variables: an outcome Y with two possible values $\{1, 2\}$ (where 2 codes a success), and covariates $X \in \{1, \dots, I\}$ and $Z \in \{1, \dots, J\}$. The data from n independent observations of (X, Z, Y) can be summarized as a 3-way contingency table of the size $I \times J \times 2$ with observed counts n_{ijk} of the combinations $X = i$, $Z = j$ and

$Y = k$, expected counts m_{ijk} , and cell probabilities π_{ijk} .

We are interested in estimating the conditional probabilities p_{ij} of success given the covariates, that is $p_{ij} = P(Y = 2 | X = i, Z = j) = m_{ij2}/m_{ij+} = \pi_{ij2}/\pi_{ij+}$. The problem can be addressed either by logistic regression or by a loglinear model. We will show that, with an appropriately selected loglinear model, the results from the two approaches are equivalent.

Let the correct logistic model be

$$\log \frac{p_{ij}}{1 - p_{ij}} = \gamma^0 + \gamma_i^X + \gamma_j^Z, \quad (3.19)$$

where $\gamma_1^X = \gamma_1^Z = 0$. The left-hand side of the model can be rewritten as

$$\log \frac{p_{ij}}{1 - p_{ij}} = \log \frac{m_{ij2}}{m_{ij1}} = \log \frac{\pi_{ij2}}{\pi_{ij1}}.$$

The regression parameters of the logistic model have the following interpretation:

$$\begin{aligned} e^{\gamma_i^X} &= \frac{P(Y = 2 | X = i, Z = j)}{P(Y = 1 | X = i, Z = j)} \frac{P(Y = 1 | X = 1, Z = j)}{P(Y = 2 | X = 1, Z = j)} \\ &= \frac{p_{ij}}{1 - p_{ij}} \frac{1 - p_{1j}}{p_{1j}} = \frac{\pi_{ij2}\pi_{1j1}}{\pi_{ij1}\pi_{1j2}} = \theta_{i2(j)}^{XY}, \end{aligned}$$

which is the conditional odds ratio for the association between X and Y given $Z = j$ (and it does not depend on j). Next,

$$\begin{aligned} e^{\gamma_j^Z} &= \frac{P(Y = 2 | X = i, Z = j)}{P(Y = 1 | X = i, Z = j)} \frac{P(Y = 1 | X = i, Z = 1)}{P(Y = 2 | X = i, Z = 1)} \\ &= \frac{p_{ij}}{1 - p_{ij}} \frac{1 - p_{i1}}{p_{i1}} = \frac{\pi_{ij2}\pi_{i11}}{\pi_{ij1}\pi_{i12}} = \theta_{j2(i)}^{ZY}, \end{aligned}$$

which is the conditional odds ratio for the association between Z and Y given $X = i$ (and it does not depend on i).

Now consider the loglinear model with all two-way interactions, that is (XY, ZY, XZ) :

$$\log m_{ijk} = \alpha + \beta_i^X + \beta_j^Z + \beta_k^Y + \beta_{ik}^{XY} + \beta_{jk}^{ZY} + \beta_{ij}^{XZ} \quad (3.20)$$

with the usual constraints on the parameters. In this model,

$$\log \frac{p_{ij}}{1 - p_{ij}} = \log \frac{m_{ij2}}{m_{ij1}} = \beta_2^Y + \beta_{i2}^{XY} + \beta_{j2}^{ZY}.$$

Thus, the loglinear model (3.20) induces the same structure on the conditional log odds of success as the logistic model (3.19). Clearly, $\beta_{i2}^{XY} = \gamma_i^X$, $\beta_{j2}^{ZY} = \gamma_j^Z$, and $\beta_2^Y = \gamma^0$. Thus, the logistic model estimates a subset of the parameters of the loglinear model that

determine the associations between Y and (X, Z) . The results about these associations (parameter estimates, hypotheses tests) are the same no matter if they were obtained from the logistic model (3.19) or the loglinear model (3.20).

The only difference between the logistic model (3.19) and the loglinear model (3.20) is that the loglinear model also estimates the associations between X and Z . These associations are not estimated by the logistic model (they are not of interest).

One can easily generalize this observation about equivalence between loglinear and logistic models to arbitrary categorical covariates.

Note. The logistic model $Y \sim M$ is equivalent to the loglinear model $(MY, \mathfrak{I}(M))$, where MY includes all the interactions between the terms in M with Y and $\mathfrak{I}(M)$ is the most general interaction between all the terms in M .

4 Extensions of Generalized Linear Models

4.1 Quasi-likelihood and Overdispersion

4.1.1 Overdispersion in binomial data

Let $Y_1, \dots, Y_n \sim \text{Bi}(m, \pi_0)$ be independent observations with moments $\mathbf{E}Y_i = m\pi_0$ and $\text{var} Y_i = m\pi_0(1 - \pi_0)$. Since each of them can be decomposed into iid alternative variables, methods for exponential family distributions apply to this case.

Now suppose that π_0 is not a constant but a random variable with mean π_0 : the i th binomial variable has the distribution $\text{Bi}(m, \pi_i)$ (given π_i), where π_1, \dots, π_n is a random sample from a distribution with a density $g(\pi)$, mean $\mathbf{E} \pi_i = \pi_0$ and variance $\text{var} \pi_i = \sigma_\pi^2$. Then $\mathbf{E}Y_i = m\pi_0$ but $\text{var} Y_i = m\pi_0(1 - \pi_0) + m(m - 1)\sigma_\pi^2$. The variability in Y_i is larger than for the binomial distribution unless $m = 1$. This effect is called *overdispersion**.

Let $m > 1$ and consider a specific distribution for π_i , for example $\pi_i \sim \text{B}(\alpha, \beta)$. Then $\pi_0 = \alpha/(\alpha + \beta)$ and $\sigma_\pi^2 = \pi_0(1 - \pi_0)/(\alpha + \beta + 1)$. The distribution of Y_i is

$$\mathbf{P}[Y_i = j] = \binom{m}{j} \frac{B(\alpha + j, \beta + m - j)}{B(\alpha, \beta)}, \quad j = 0, \dots, m.$$

This is called a *beta-binomial distribution*. Its variance is

$$\varphi m\pi_0(1 - \pi_0), \quad \text{where} \quad \varphi = 1 + \frac{m - 1}{\alpha + \beta + 1} > 1$$

plays the role of a dispersion parameter.

This is not a distribution from the exponential family but we would like to extend the theory of GLM to responses following distributions of this type.

4.1.2 Overdispersion in Poisson data

Let $Y_1, \dots, Y_n \sim \text{Po}(\lambda_0)$ be independent observations with moments $\mathbf{E}Y_i = \lambda_0$ and $\text{var} Y_i = \lambda_0$. This is a distribution from the exponential family.

Let the Poisson parameters for the observations be random variables with mean λ_0 rather than a constant: the i th Poisson variable has the distribution $\text{Po}(\lambda_i)$ (given

* *Česky nadměrná disperse*

λ_i), where $\lambda_1, \dots, \lambda_n$ is a random sample from a distribution with a density $g(\lambda)$, mean $\mathbf{E} \lambda_i = \lambda_0$ and variance $\text{var} \lambda_i = \sigma_\lambda^2$. Then $\mathbf{E} Y_i = \lambda_0$ but $\text{var} Y_i = \lambda_0 + \sigma_\lambda^2 > \lambda_0$. The variability in Y_i is larger than for the Poisson distribution, another example of overdispersion.

Now suppose $\lambda_i \sim \Gamma(a, a\lambda_0)$. Then λ_0 is the mean of λ_i and $\sigma_\lambda^2 = \lambda_0/a$. The distribution of Y_i is

$$\mathbf{P}[Y_i = j] = \frac{1}{jB(j, a\lambda_0)} \left(\frac{a}{a+1} \right)^{a\lambda_0} \frac{1}{(a+1)^j}, \quad j = 0, 1, 2, \dots$$

This is called a *Poisson-gamma distribution*. Its variance is $\varphi\lambda_0$, where $\varphi = 1 + 1/a > 1$ plays the role of a dispersion parameter. Special cases of this distribution are negative binomial $\text{NB}(m, p)$ (for $a\lambda_0 = m \in \mathbf{N}$, $a/(a+1) = p$) and geometric $\text{Geo}(p)$ (for $a\lambda_0 = 1$, $a/(a+1) = p$).

Again, this distribution does not belong to the exponential family. The results for the GLM need to be extended to responses with distributions of this type.

4.1.3 Quasi-likelihood

Consider n independent copies of random vectors (Y_i, \mathbf{X}_i) , $i = 1, \dots, n$, where $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^\top$ are the covariates and Y_i is the response.

Assumptions.

1. Y_1, \dots, Y_n are independent
2. The mean $\mu_i = \mathbf{E} Y_i$ satisfies the identity $g(\mu_i) = \eta_i$, where $\eta_i = \mathbf{X}_i^\top \boldsymbol{\beta}_0$ is the linear predictor and g is a known strictly monotone, twice continuously differentiable link function.
3. The variance $\text{var} Y_i$ satisfies the identity $\text{var} Y_i = \varphi V(\mu_i)$, where $\varphi > 0$ is a dispersion parameter and V is a known positive continuously differentiable variance function.

Note. The form of the distribution of Y_i is not specified, it does not have to belong to the exponential family. Instead, a variance function describing the relationship of the mean to the variance is supplied.

Note. The original GLM was a parametric model. This is a semi-parametric model: the form of the distribution of Y_i is not specified, we only specify conditions on the first two moments of Y_i .

Because this is not a parametric model maximum likelihood estimator cannot be used. The regression parameters are estimated by the *maximum quasi-likelihood estimator*.

Definition 4.1. (Wedderburn 1974) The quasi-(log)likelihood $Q(\boldsymbol{\beta})$ is defined as $Q(\boldsymbol{\beta}) = \sum_{i=1}^n Q_i(\boldsymbol{\beta})$, where

$$Q_i(\boldsymbol{\beta}) = \int_{Y_i}^{\mu_i} \frac{Y_i - t}{\varphi V(t)} dt.$$

The maximum quasi-likelihood estimator $\hat{\boldsymbol{\beta}}_n$ is the point that maximizes the quasi-likelihood*.

The maximum quasi-likelihood estimator solves the system of equations $\mathbf{U}_n(\hat{\boldsymbol{\beta}}_n) = \mathbf{0}$, where

$$\mathbf{U}_n(\boldsymbol{\beta}) = \sum_{i=1}^n \mathbf{U}_i(\boldsymbol{\beta})$$

is the quasi-score† with the terms

$$\mathbf{U}_i(\boldsymbol{\beta}) = \frac{\partial Q_i(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \frac{Y_i - \mu_i}{\varphi V(\mu_i)} \frac{\partial \mu_i}{\partial \boldsymbol{\beta}}.$$

The quasi-score has exactly the same form as the score in the GLM, that is, it can be written as a sum of the terms

$$\frac{1}{\varphi} w(\mu_i) g'(\mu_i) (Y_i - \mu_i) \mathbf{X}_i.$$

This fact was the primary motivation to introduce the quasi-likelihood in the form specified in Definition 4.1. The following lemma says that the quasi-score has exactly the same properties as the GLM score.

Lemma 4.1. (Wedderburn 1974)

- (i) $\mathbf{U}_i(\boldsymbol{\beta})$, $i = 1, \dots, n$, are iid random vectors.
- (ii) If $\boldsymbol{\beta}_0$ is the true parameter then $\mathbb{E} \mathbf{U}_i(\boldsymbol{\beta}_0) = \mathbf{0}$.
- (iii) If $\boldsymbol{\beta}_0$ is the true parameter then $\text{var} \mathbf{U}_i(\boldsymbol{\beta}_0) = -\mathbb{E} \frac{\partial}{\partial \boldsymbol{\beta}^\top} \mathbf{U}_i(\boldsymbol{\beta}_0) = I(\boldsymbol{\beta}_0)$, where $I(\boldsymbol{\beta}_0)$ is defined by (2.8).

Proposition 4.2. There exists a sequence $\hat{\boldsymbol{\beta}}_n$ of solutions to the quasi-likelihood equations such that $\hat{\boldsymbol{\beta}}_n \xrightarrow{\text{P}} \boldsymbol{\beta}_0$.

The proposition claims consistency as long as the solutions exist and are unique. The key for the validity of the proposition is Lemma 4.1 (ii). Proposition 4.2 together with Lemma 4.1 justify the validity of two of the key asymptotic results that hold in the GLM. The results that we are getting in this section generalize Proposition 1.2.

* Český kvazivěrohodnost † Český kvaziskóre

Theorem 4.3.

- (i) $\frac{1}{\sqrt{n}}\mathbf{U}_n(\boldsymbol{\beta}_0) \xrightarrow{D} \mathbf{N}_p(\mathbf{0}, I(\boldsymbol{\beta}_0)),$
- (ii) $\sqrt{n}(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) \xrightarrow{D} \mathbf{N}_p(\mathbf{0}, I^{-1}(\boldsymbol{\beta}_0)).$

Thus, Wald tests and score tests (and confidence intervals) derived for the GLM also hold when quasi-likelihood is used instead of full likelihood. On the other hand, likelihood ratio (and hence deviance) tests do not work.

The dispersion parameter φ can be estimated by the method of moments based on Pearson X^2 statistic as described in Section 2.5.

The theory of GLM (except deviance tests) holds even for distributions that are not of exponential type as long as the data are independent and the variance function is correctly specified. Thus, we can fit GLM to distributions such as beta-binomial, negative binomial or geometric (and many more).

4.2 Sandwich Variance Estimation in the GLM

4.2.1 Behavior of the MLE under a misspecified model

Let X_1, \dots, X_n be iid random variables (vectors) on the space $(\mathcal{X}, \mathcal{A})$ with distribution P and density p with respect to a σ -finite measure μ .

Consider the model $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ on the space $(\mathcal{X}, \mathcal{A})$ with densities p_θ with respect to μ . Let $\Theta \subseteq \mathbb{R}^d$. Suppose the model \mathcal{P} is regular. We do not assume that $P \in \mathcal{P}$. There need not exist any $\theta \in \Theta$ such that $P = P_\theta$.

We use the data $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} P$ and the model \mathcal{P} to estimate θ by the method of maximum likelihood. This section summarizes the results obtained by [White \(1982\)](#).

Define

$$\begin{aligned} \ell_n(\boldsymbol{\theta}) &= \sum_{i=1}^n \log p_{\boldsymbol{\theta}}(X_i), \\ \widehat{\boldsymbol{\theta}}_n &= \arg \max_{\boldsymbol{\theta} \in \Theta} \ell_n(\boldsymbol{\theta}), \\ \mathbf{U}_i(\boldsymbol{\theta}) &= \frac{\partial}{\partial \boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(X_i), \\ \mathbf{U}_n(\boldsymbol{\theta}) &= \sum_{i=1}^n \mathbf{U}_i(\boldsymbol{\theta}). \end{aligned}$$

The estimator $\widehat{\boldsymbol{\theta}}_n$ solves the system of equations $\mathbf{U}_n(\widehat{\boldsymbol{\theta}}_n) = \mathbf{0}$. [White \(1982\)](#) calls it the *quasi-maximum likelihood estimator* but we will prefer the term *pseudo-maximum likelihood estimator*. What does this estimator estimate when the model does not hold?

Theorem 4.4. (White 1982) $\hat{\boldsymbol{\theta}}_n \xrightarrow{P} \boldsymbol{\theta}_0$ as $n \rightarrow \infty$, where

$$\boldsymbol{\theta}_0 = \arg \min_{\boldsymbol{\theta} \in \Theta} K(P, P_{\boldsymbol{\theta}}) = \arg \max_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_P \log p_{\boldsymbol{\theta}}(X_i)$$

and

$$K(P, P_{\boldsymbol{\theta}}) = \mathbb{E}_P \log \frac{p(X_i)}{p_{\boldsymbol{\theta}}(X_i)} \geq 0.$$

The pseudo-MLE converges to the point $\boldsymbol{\theta}_0$ that minimizes the Kullback-Leibler distance between all the distributions belonging to the model and the true distribution of the data.

The pseudo-score statistic is asymptotically normal when evaluated at $\boldsymbol{\theta}_0$ and the estimator is also asymptotically normal.

Theorem 4.5. (White 1982)

(i) The probability limit $\boldsymbol{\theta}_0$ of $\hat{\boldsymbol{\theta}}_n$ satisfies $\mathbb{E}_P \mathbf{U}_n(\boldsymbol{\theta}_0) = \mathbf{0}$.

(ii)

$$\frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\theta}_0) \xrightarrow{D} \mathbf{N}_d(\mathbf{0}, \Sigma),$$

where

$$\Sigma = \text{var}_P \frac{\partial}{\partial \boldsymbol{\theta}} \log p_{\boldsymbol{\theta}_0}(X_i).$$

(iii)

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow{D} \mathbf{N}_d(\mathbf{0}, \mathbb{D}^{-1} \Sigma \mathbb{D}^{-1}),$$

where

$$\mathbb{D} = -\mathbb{E}_P \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \log p_{\boldsymbol{\theta}_0}(X_i).$$

The asymptotic variance of $\hat{\boldsymbol{\theta}}_n$ has the sandwich form. The matrix \mathbb{D} plays the role of an information matrix. However, \mathbb{D} is not equal to the asymptotic variance Σ of the pseudo-score statistic.

Theorem 4.6. The asymptotic variance matrix $\mathbb{D}^{-1} \Sigma \mathbb{D}^{-1}$ of $\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$ can be consistently estimated by $\hat{\mathbb{D}}^{-1} \hat{\Sigma} \hat{\mathbb{D}}^{-1}$, where

$$\hat{\mathbb{D}} = -\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \log p_{\hat{\boldsymbol{\theta}}_n}(X_i)$$

and

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \mathbf{U}_i(\hat{\boldsymbol{\theta}}_n)^{\otimes 2}.$$

This is the *sandwich estimator** of the asymptotic variance. Compare this to Proposition 1.3 and its use in linear regression.

* Český sendvičový odhad

Note.

- If the model holds, that is $\exists \boldsymbol{\theta}_0 \in \Theta$ such that $P = P_{\boldsymbol{\theta}_0}$, then $\widehat{\boldsymbol{\theta}}_n \xrightarrow{P} \boldsymbol{\theta}_0$ and $\mathbb{D} \equiv \Sigma$.
- The sandwich estimator $\widehat{\mathbb{D}}^{-1} \widehat{\Sigma} \widehat{\mathbb{D}}^{-1}$ of the asymptotic variance tends to underestimate the true variance $\mathbb{D}^{-1} \Sigma \mathbb{D}^{-1}$ unless n is very large. Various modifications have been proposed to reduce the small sample bias of the sandwich estimator.
- This theory is especially useful when at least some components of $\boldsymbol{\theta}_0$ are equal to the true parameters that we wish to estimate.

4.2.2 Applications to the GLM

Consider n independent copies of random vectors (Y_i, \mathbf{X}_i) , $i = 1, \dots, n$, where $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^\top$ are the covariates and Y_i is the response.

Assumptions.

1. Y_1, \dots, Y_n are independent
2. The mean $\mu_i = \mathbf{E} Y_i$ satisfies the identity $g(\mu_i) = \eta_i$, where $\eta_i = \mathbf{X}_i^\top \boldsymbol{\beta}_0$ is the linear predictor and g is a known strictly monotone, twice continuously differentiable link function.

Note.

- We will choose a *working variance function** $V(\mu_i)$ and use it in the estimation of $\boldsymbol{\beta}_0$ but we will not assume that this variance function is correct.
- No assumptions are made about the form of the density of Y_i .

The estimation of $\boldsymbol{\beta}_0$ proceeds as if Y_i had a distribution from the exponential family with mean $\mu_i = \mathbf{E} Y_i = g^{-1}(\eta_i)$ and variance $\text{var} Y_i = \varphi V(\mu_i)$. The *pseudo-score function*† is the same as in the GLM,

$$\mathbf{U}_i(\boldsymbol{\beta}) = \frac{Y_i - \mu_i}{\varphi V(\mu_i)} \frac{\partial \mu_i}{\partial \boldsymbol{\beta}} = \frac{1}{\varphi} w(\mu_i) g'(\mu_i) (Y_i - \mu_i) \mathbf{X}_i.$$

The estimator $\widehat{\boldsymbol{\beta}}_n$ is the solution to the system of pseudo-score equations

$$\mathbf{U}_n(\widehat{\boldsymbol{\beta}}_n) = \sum_{i=1}^n w(\widehat{\mu}_i) g'(\widehat{\mu}_i) (Y_i - \widehat{\mu}_i) \mathbf{X}_i = \mathbf{0},$$

where $\widehat{\mu}_i = g^{-1}(\mathbf{X}_i^\top \widehat{\boldsymbol{\beta}}_n)$. This system can be solved by the IWLS algorithm.

Lemma 4.7.

- (i) $\mathbf{U}_i(\boldsymbol{\beta})$, $i = 1, \dots, n$, are iid random vectors.

* Český pracovní rozptylová funkce † Český pseudoskórová funkce

(ii) If β_0 is the true parameter then $\mathbb{E} U_i(\beta_0) = \mathbf{0}$.

We use the theory from the previous section to derive the probability limit and the asymptotic distribution of $\hat{\beta}_n$. Let

$$I(\beta_0) = -\mathbb{E} \frac{\partial}{\partial \beta^\top} U_i(\beta_0) = \frac{1}{\varphi} \mathbb{E}_{\mathbf{X}} w(\mu_i) \mathbf{X}_i^{\otimes 2}$$

and

$$\Sigma = \text{var} U_i(\beta_0) = \frac{1}{\varphi^2} \mathbb{E}_{\mathbf{X}} w^2(\mu_i) [g'(\mu_i)]^2 \text{var} Y_i \mathbf{X}_i^{\otimes 2} \neq I(\beta_0).$$

It follows from Theorem 4.5 (note that $I(\beta_0)$ plays the role of the matrix \mathbb{D}) that

Theorem 4.8.

- (i) $\hat{\beta}_n$ converges in probability to β_0 .
- (ii) $\frac{1}{\sqrt{n}} U_n(\beta_0) \xrightarrow{\text{D}} \mathbf{N}_p(\mathbf{0}, \Sigma)$,
- (iii) $\sqrt{n}(\hat{\beta}_n - \beta_0) \xrightarrow{\text{D}} \mathbf{N}_p(\mathbf{0}, I^{-1}(\beta_0) \Sigma I^{-1}(\beta_0))$.

Corollary. The asymptotic variance matrix of $\sqrt{n}(\hat{\beta}_n - \beta_0)$ can be consistently estimated by $\hat{I}^{-1} \hat{\Sigma} \hat{I}^{-1}$, where \hat{I} is defined by (2.9) and

$$\hat{\Sigma} = \frac{1}{n \hat{\varphi}^2} \sum_{i=1}^n [w(\hat{\mu}_i) g'(\hat{\mu}_i) (Y_i - \hat{\mu}_i)]^2 \mathbf{X}_i^{\otimes 2}.$$

The corollary follows from Theorem 4.6. The dispersion parameter φ can be estimated by the Pearson X^2 statistic (see Section 2.5) but it is not needed to calculate either $\hat{\beta}_n$ or its estimated asymptotic variance $\hat{I}^{-1} \hat{\Sigma} \hat{I}^{-1}$.

Likelihood ratio (deviance) tests cannot be used but Wald tests and score tests based on Theorem 4.8 are available. Thus, even if the distribution of the responses is not of exponential family and the variance function is unknown, the theory of the GLM can be used for parameter estimation and asymptotic variance can be estimated by the sandwich estimator. If the working variance function $V(\mu)$ is guessed correctly then $\Sigma \approx I(\beta_0)$ and the results will be close to those obtained by the quaslikelihood approach. If the working variance function $V(\mu)$ is far from the truth the asymptotic variance will increase and estimates will be less efficient. The most serious danger of this approach is the potential underestimation of the true variance by the sandwich.

5 Generalized Estimating Equations

5.1 Group-Dependent Data

In this chapter, we extend the generalized linear model even further: we will develop a regression model that can be used for certain kinds of correlated data.

Suppose we observe K independent random vectors $\mathbf{Y}_1, \dots, \mathbf{Y}_K$, where

$$\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^\top, \quad i = 1, \dots, K.$$

So the data consists of K independent groups (subjects), which include different numbers of correlated observations. Within each group, the observations are dependent but between groups, they are independent.

We will call such data structures *group-dependent data*^{*} but they may be called by several different terms depending on the context or application field.

Clustered data are measurements collected on groups of objects such as teeth, animal litters, or siblings. The measurements in the same group do not have a well defined ordering.[†]

Repeated measures are measurements made repeatedly on the same subject so that there is ordering among them (first, second, third, ...).[‡]

Longitudinal data are repeated measures with captured time information. This means, each measurement is made at a different time, which is recorded in the data. Longitudinal data arise by observing independent short pieces of time series.[§]

Panel data is the term used for group-dependent data in econometrics.[¶]

Each observation Y_{ij} is accompanied by a vector of covariates \mathbf{X}_{ij} of the size p . We would like to describe the dependence of $\mu_{ij} = \mathbb{E}Y_{ij}$ on the covariates \mathbf{X}_{ij} by a regression model. As in the GLM, we will assume that $g(\mu_{ij}) = \mathbf{X}_{ij}^\top \boldsymbol{\beta}_0$, where g is a known strictly monotone, twice continuously differentiable link function and $\boldsymbol{\beta}_0$ is an unknown true parameter vector.

Thus, we assume that

$$\mathbb{E} \mathbf{Y}_i = \boldsymbol{\mu}_i = (\mu_{i1}, \dots, \mu_{in_i})^\top,$$

where $\mu_{ij} = g^{-1}(\mathbf{X}_{ij}^\top \boldsymbol{\beta}_0)$. Like in Section 4.2.2, We leave $\text{var} \mathbf{Y}_i$ unspecified – we do not impose any assumptions on variances and covariances of the measurements.

^{*} Česky *skupinově závislá data* [†] Česky *shluková data* [‡] Česky *opakovaná měření* [§] Česky *longitudinální data* [¶] Česky *panelová data*

5.2 Estimation of Regression Parameters by Generalized Estimating Equations

A method for estimating regression parameters under this general extension of the GLM to group-dependent data was first proposed by Liang and Zeger (1986). They called it *generalized estimating equations** [GEE]. Let

$$\mathbb{X}_i = \begin{pmatrix} \mathbf{X}_{i1}^\top \\ \vdots \\ \mathbf{X}_{in_i}^\top \end{pmatrix}_{n_i \times p} \quad \text{and} \quad \left(\frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} \right)_{p \times n_i} = \mathbb{X}_i^\top \begin{pmatrix} g'(\mu_{i1}) & \dots & 0 \\ 0 & \ddots & 0 \\ 0 & \dots & g'(\mu_{in_i}) \end{pmatrix}^{-1}.$$

Recall that the score function in the GLM (with $n_i = 1$ for all i) can be written as

$$\mathbf{U}_i(\boldsymbol{\beta}) = \frac{\partial \mu_i}{\partial \boldsymbol{\beta}} \frac{1}{\varphi V(\mu_i)} (Y_i - \mu_i).$$

When written in this way, the score can be easily generalized to multivariate \mathbf{Y}_i . We define a pseudo-score function

$$\mathbf{U}_i(\boldsymbol{\beta}) = \left(\frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} \right) \mathbb{Q}_i^{-1}(\boldsymbol{\mu}_i) (\mathbf{Y}_i - \boldsymbol{\mu}_i),$$

where

$$\mathbb{Q}_i(\boldsymbol{\mu}_i) = \varphi \mathbb{V}_i^{1/2}(\boldsymbol{\mu}_i) \mathbb{R}_i \mathbb{V}_i^{1/2}(\boldsymbol{\mu}_i)$$

represents our guess about $\text{var } \mathbf{Y}_i$. The diagonal $n_i \times n_i$ matrix \mathbb{V}_i provides working variances of the observations, $V(\mu_{i1}), \dots, V(\mu_{in_i})$, on its diagonal. The $n_i \times n_i$ matrix \mathbb{R}_i is a *working correlation matrix*[†], our guess about $\text{cor}(Y_{ij}, Y_{ik})$.

Neither the variance function $V(\cdot)$ nor the correlation matrix \mathbb{R}_i is assumed to be correct. If we knew $\text{var } Y_{ij}$ and $\text{cor } \mathbf{Y}_i$ we would use them in place of $\varphi V(\mu_{ij})$ and \mathbb{R}_i . Since we do not know them, we just use our best guess.

Because $\mathbb{E} \mathbf{Y}_i = \boldsymbol{\mu}_i$, we easily get the moments of $\mathbf{U}_i(\boldsymbol{\beta}_0)$:

$$\mathbb{E} \mathbf{U}_i(\boldsymbol{\beta}_0) = \mathbf{0} \quad \text{and} \quad \text{var } \mathbf{U}_i(\boldsymbol{\beta}_0) \equiv \Sigma = \mathbb{E} \left(\frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} \right) \mathbb{Q}_i^{-1}(\boldsymbol{\mu}_i) \text{var } \mathbf{Y}_i \mathbb{Q}_i^{-1}(\boldsymbol{\mu}_i) \left(\frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} \right)^\top.$$

The estimator $\widehat{\boldsymbol{\beta}}_K$ is defined as the solution to the system of equations

$$\mathbf{U}_K(\widehat{\boldsymbol{\beta}}_K) \equiv \sum_{i=1}^K \mathbf{U}_i(\widehat{\boldsymbol{\beta}}_K) = \mathbf{0}. \quad (5.1)$$

This is called *the GEE estimator*. It is another special case of the methods of Section 4.2.1. Heuristically, the properties of $\widehat{\boldsymbol{\beta}}_K$ follow from Theorem 4.5, though a rigorous proof of consistency would require more work.

* Český zobecněné odhadovací rovnice † Český pracovní korelační matice

Let

$$\mathbb{D} = -\mathbb{E} \frac{\partial}{\partial \boldsymbol{\beta}^\top} \mathbf{U}_i(\boldsymbol{\beta}_0) = \mathbb{E} \left(\frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} \right) \mathbb{Q}_i^{-1}(\boldsymbol{\mu}_i) \left(\frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} \right)^\top.$$

Proposition 5.1. (*Liang and Zeger 1986*) As $K \rightarrow \infty$,

- (i) $\hat{\boldsymbol{\beta}}_K \xrightarrow{\text{P}} \boldsymbol{\beta}_0$,
- (ii) $\frac{1}{\sqrt{K}} \mathbf{U}_K(\boldsymbol{\beta}_0) \xrightarrow{\text{D}} \mathbf{N}_p(\mathbf{0}, \Sigma)$,
- (iii) $\sqrt{K}(\hat{\boldsymbol{\beta}}_K - \boldsymbol{\beta}_0) \xrightarrow{\text{D}} \mathbf{N}_p(\mathbf{0}, \mathbb{D}^{-1} \Sigma \mathbb{D}^{-1})$.

Notice that the asymptotics requires that the number K of independent groups tends to infinity. The number of observations within the groups is irrelevant. By Theorem 4.6, the asymptotic variance of $\mathbb{D}^{-1} \Sigma \mathbb{D}^{-1}$ can be consistently estimated by the sandwich $\hat{\mathbb{D}}^{-1} \hat{\Sigma} \hat{\mathbb{D}}^{-1}$, where

$$\hat{\mathbb{D}} = \frac{1}{K} \sum_{i=1}^K \left(\frac{\partial \hat{\boldsymbol{\mu}}_i}{\partial \boldsymbol{\beta}} \right) \mathbb{Q}_i^{-1}(\hat{\boldsymbol{\mu}}_i) \left(\frac{\partial \hat{\boldsymbol{\mu}}_i}{\partial \boldsymbol{\beta}} \right)^\top$$

and

$$\hat{\Sigma} = \frac{1}{K} \sum_{i=1}^K \mathbf{U}_i(\hat{\boldsymbol{\beta}}_K)^{\otimes 2}.$$

The system of equations (5.1) can be solved by a modified IWLS method. The solver iterates

$$\hat{\boldsymbol{\beta}} = \left[\sum_{i=1}^K \left(\frac{\partial \hat{\boldsymbol{\mu}}_i}{\partial \boldsymbol{\beta}} \right) \mathbb{Q}_i^{-1}(\hat{\boldsymbol{\mu}}_i) \left(\frac{\partial \hat{\boldsymbol{\mu}}_i}{\partial \boldsymbol{\beta}} \right)^\top \right]^{-1} \left[\sum_{i=1}^K \left(\frac{\partial \hat{\boldsymbol{\mu}}_i}{\partial \boldsymbol{\beta}} \right) \mathbb{Q}_i^{-1}(\hat{\boldsymbol{\mu}}_i) \hat{\mathbf{Z}}_i \right],$$

where

$$\hat{\mathbf{Z}}_i = (\hat{Z}_{i1}, \dots, \hat{Z}_{in_i})^\top \quad \text{and} \quad \hat{Z}_{ij} = \frac{\hat{\eta}_{ij} + (Y_{ij} - \hat{\mu}_{ij})g'(\hat{\mu}_{ij})}{g'(\hat{\mu}_{ij})}.$$

The dispersion parameter φ can be estimated by

$$\hat{\varphi} = \frac{1}{n-p} \sum_{i=1}^K \sum_{j=1}^{n_i} \frac{(Y_{ij} - \hat{\mu}_{ij})^2}{V(\hat{\mu}_{ij})},$$

which is based on a modified Pearson X^2 statistic (ignoring the correlations). Here, $n = \sum_{i=1}^K n_i$ is the total sample size. This estimator is consistent if the variance function V is correctly specified, but it is not efficient. Estimated φ is not needed to calculate either $\hat{\boldsymbol{\beta}}_n$ or its estimated asymptotic variance $\hat{\mathbb{D}}^{-1} \hat{\Sigma} \hat{\mathbb{D}}^{-1}$.

5.3 Correlation Structures

How should we choose the working covariance matrix $\mathbb{Q}_i(\boldsymbol{\mu}_i) = \varphi \mathbb{V}_i^{1/2}(\boldsymbol{\mu}_i) \mathbb{R}_i \mathbb{V}_i^{1/2}(\boldsymbol{\mu}_i)$? The variance function $V(\cdot)$ expresses our belief about the dependence of the variance on

the mean. It is more complicated to choose a suitable working correlation matrix \mathbb{R}_i for each of the independent groups. The current section is devoted to this problem.

5.3.1 Working independence

We take $\mathbb{R}_i = \mathbb{I}_{n_i}$ (identity matrix). Then $\mathbb{Q}_i(\boldsymbol{\mu}_i) = \varphi \mathbb{V}_i(\boldsymbol{\mu}_i)$ is diagonal and $\boldsymbol{\beta}_0$ is estimated as if the data were independent, by fitting the standard IWLS. After the estimates are obtained under the independence assumption (they are still consistent), their variance is adjusted by the sandwich to take into account correlations and also perhaps the wrong choice of $V(\cdot)$. This is the easiest way to fit GEE. The estimates will be consistent but inefficient if the correlations are strong.

5.3.2 Parametrized correlations

The other option is to introduce some non-independent correlation structure and parametrize it by an m -dimensional parameter vector $\boldsymbol{\alpha} \in \mathbb{R}^m$. We take $\mathbb{R}_i = \mathbb{R}_i(\boldsymbol{\alpha})$,

$$\mathbb{Q}_i(\boldsymbol{\mu}_i) \equiv Q_i(\boldsymbol{\mu}_i, \boldsymbol{\alpha}) = \varphi \mathbb{V}_i^{1/2}(\boldsymbol{\mu}_i) \mathbb{R}_i(\boldsymbol{\alpha}) \mathbb{V}_i^{1/2}(\boldsymbol{\mu}_i)$$

and

$$\widehat{\mathbb{Q}}_i(\boldsymbol{\mu}_i) \equiv Q_i(\boldsymbol{\mu}_i, \widehat{\boldsymbol{\alpha}}) = \widehat{\varphi} \mathbb{V}_i^{1/2}(\boldsymbol{\mu}_i) \mathbb{R}_i(\widehat{\boldsymbol{\alpha}}) \mathbb{V}_i^{1/2}(\boldsymbol{\mu}_i),$$

where $\widehat{\boldsymbol{\alpha}}$ is a \sqrt{K} -consistent estimator of $\boldsymbol{\alpha}$, for example some moment estimator.

The score is modified as follows

$$\widehat{\mathbf{U}}_i(\boldsymbol{\beta}) = \left(\frac{\partial \boldsymbol{\mu}_i}{\partial \boldsymbol{\beta}} \right) \widehat{\mathbb{Q}}_i^{-1}(\boldsymbol{\mu}_i) (\mathbf{Y}_i - \boldsymbol{\mu}_i),$$

but these vectors are no longer independent for $i = 1, \dots, n$. One needs to show that

$$\frac{1}{\sqrt{K}} \sum_{i=1}^K \widehat{\mathbf{U}}_i(\boldsymbol{\beta}_0) = \frac{1}{\sqrt{K}} \sum_{i=1}^K \mathbf{U}_i(\boldsymbol{\beta}_0) + o_P(1).$$

(see [Liang and Zeger 1986](#)). The estimator $\widehat{\boldsymbol{\beta}}$ solves $\sum_{i=1}^K \widehat{\mathbf{U}}_i(\widehat{\boldsymbol{\beta}}) = \mathbf{0}$ and Proposition 5.1 still holds.

Here is a general strategy how to estimate parametrized correlations by the method of moments:

1. Estimate $\boldsymbol{\beta}$ under working independence (take $\mathbb{R}_i = \mathbb{I}_{n_i}$).
2. Calculate Pearson residuals

$$r_{ij}^P = \frac{Y_{ij} - \widehat{\mu}_{ij}}{\sqrt{V(\widehat{\mu}_{ij})}}.$$

3. If the mean structure and variance function are correct we have

$$\mathbf{E} r_{ij}^P \approx 0, \quad \text{var } r_{ij}^P \approx \varphi, \quad \mathbf{E} r_{ij}^P r_{ik}^P \approx \varphi R_{ijk}(\boldsymbol{\alpha}).$$

We can use moment estimators of $\boldsymbol{\alpha}$ based on the products $r_{ij}^P r_{ik}^P$ of Pearson residuals from the same group. The next section shows examples of such estimators for selected correlation structures.

1-band correlation

Let

$$\text{cor}(Y_{ij}, Y_{ik}) = \begin{cases} 1 & \text{if } j = k, \\ \alpha & \text{if } |j - k| = 1, \\ 0 & \text{if } |j - k| > 1. \end{cases}$$

The parameter α is consistently estimated by

$$\hat{\alpha} = \frac{1}{\hat{\varphi}} \frac{1}{N - K - p} \sum_{i=1}^K \sum_{j=1}^{n_i-1} r_{ij}^P r_{i,j+1}^P.$$

m -band correlation

Take some positive whole number m . Let

$$\text{cor}(Y_{ij}, Y_{ik}) = \begin{cases} 1 & \text{if } j = k, \\ \alpha_l & \text{if } |j - k| = l, \quad l = 1, \dots, m, \\ 0 & \text{if } |j - k| > m. \end{cases}$$

For $l \in \{1, \dots, m\}$, the parameter α_l is consistently estimated by

$$\hat{\alpha}_l = \frac{1}{\hat{\varphi}} \frac{1}{N - Kl - p} \sum_{i=1}^K \sum_{j=1}^{n_i-l} r_{ij}^P r_{i,j+l}^P.$$

Note that m should not be too large; especially it should not exceed too many n_i 's.

Exchangeable correlation

Let

$$\text{cor}(Y_{ij}, Y_{ik}) = \begin{cases} 1 & \text{if } j = k, \\ \alpha & \text{if } j \neq k. \end{cases}$$

The parameter α is consistently estimated by

$$\hat{\alpha} = \frac{1}{\hat{\varphi}} \frac{1}{\sum_{i=1}^K \binom{N_i}{2} - p} \sum_{i=1}^K \sum_{j=1}^{n_i} \sum_{k=j+1}^{n_i} r_{ij}^P r_{ik}^P.$$

AR(1) correlation

Let Y_{i1}, \dots, Y_{in_i} form an AR(1) series. Then

$$\text{cor}(Y_{ij}, Y_{ik}) = \alpha^{|j-k|}.$$

Since

$$\begin{aligned} \mathbf{E} r_{ij}^P r_{ik}^P &\approx \varphi \alpha^{|j-k|} \\ \text{and } \log \mathbf{E} r_{ij}^P r_{ik}^P &\approx \log \varphi + |j-k| \log \alpha, \end{aligned}$$

we can estimate $\log \alpha$ as the slope in the linear regression problem with $\log r_{ij}^P r_{ik}^P$ as the response and $|j-k|$ as the covariate.

Other correlation structures

See methods for estimating correlations from time series data.

5.3.3 Joint estimation of mean and correlation structures

[Prentice \(1988\)](#) and [Prentice and Zhao \(1991\)](#) proposed to estimate the correlation parameters by introducing another set of estimating equations for α and solving them jointly with the generalized estimating equations (5.1) for β . [Yan and Fine \(2004\)](#) extended this idea even further: they proposed a third set of estimating equations for φ and allow φ to depend on the covariates and thus vary between subjects.

For details, see the references.

5.3.4 Summary of GEE methods

1. GEE works for regression analysis of data with K independent groups, which are correlated within each group. The number K of independent groups must be large enough for the asymptotics to work.
2. It is not necessary to correctly specify the distribution of the response, the variance of the response, or the correlations within each group. If the variance and the correlations are seriously misspecified, the variance of the estimators increases but the estimators are still consistent and asymptotically normal.
3. The parameters β have population-averaged interpretation, not subject-specific interpretation (see the discussion in Chapter 7).

6 Linear Mixed Effects Models

6.1 Introduction

6.1.1 One-way ANOVA

Let $Y_{ij} \sim N(\mu_i, \sigma_e^2)$ be independent. The index $i = 1, \dots, I$ denotes a subject. Subjects represent the levels of factor A that may potentially affect the mean of the outcome. The index $j = 1, \dots, n_i$ denotes observations within subject. The total number of observations is $n = \sum_{i=1}^I n_i$. The null hypothesis of interest is

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_I,$$

that is, the means $E Y_{ij}$ are the same for all subjects (levels of factor A).

This is the setup of the *one-way ANOVA problem** familiar from an introductory statistics class. The model can be written in the form

$$Y_{ij} = \mu + \beta_i + \varepsilon_{ij},$$

where β_1, \dots, β_I are constants subject to the linear constraint $\sum_{i=1}^I \beta_i = 0$, μ is the overall mean, and ε_{ij} are random errors (expressing the unexplained variability in Y_{ij}) satisfying $E \varepsilon_{ij} = 0$, $\text{var } \varepsilon_{ij} = \sigma_e^2$.

There are $I + 1$ unknown parameters in this model, in particular $\mu, \beta_1, \dots, \beta_{I-1}, \sigma_e^2$.

The hypothesis test in one-way ANOVA can be presented in the form of ANOVA table (Table 6.1). The second column (SS = sums of squares) contains sums of squares

* Český analýza rozptylu – jednoduché třídění

Table 6.1: One way ANOVA table with fixed effects

Source	SS	df	MS	EMS	F
Subject (A)	SS_A	$I - 1$	$MS_A = \frac{SS_A}{I - 1}$	$\sigma_e^2 + Q_\beta$	$F_A = \frac{MS_A}{MS_e}$
Residual	SS_e	$n - I$	$MS_e = \frac{SS_e}{n - I}$	σ_e^2	
Total	SS_T	$n - 1$			

for subjects

$$SS_A = \sum_{i=1}^I n_i (\bar{Y}_i - \bar{Y}_{..})^2,$$

where $\bar{Y}_i = n_i^{-1} \sum_{j=1}^{n_i} Y_{ij}$ and $\bar{Y}_{..} = n^{-1} \sum_{i=1}^I \sum_{j=1}^{n_i} Y_{ij}$, and residual sums of squares

$$SS_e = \sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2.$$

It can be shown that

$$SS_A = \mathbf{Z}^\top \mathbb{A} \mathbf{Z},$$

where

$$\mathbf{Z} = \begin{pmatrix} \bar{Y}_1 \\ \vdots \\ \bar{Y}_I \end{pmatrix} - \mu \mathbf{1}_I \quad \text{and} \quad \mathbb{A} = \text{diag}(n_1, \dots, n_I) - \frac{1}{n} \mathbf{n} \mathbf{n}^\otimes,$$

$\mathbf{n} = (n_1, \dots, n_I)^\top$ and $\mathbf{1}_I$ is a column vector of 1's of length I .

Also,

$$SS_e = \sum_{i=1}^I \mathbf{z}_i^\top \mathbb{A}_i \mathbf{z}_i,$$

where

$$\mathbf{z}_i = \begin{pmatrix} Y_{i1} \\ \vdots \\ Y_{in_i} \end{pmatrix} - \mu_i \mathbf{1}_{n_i} \quad \text{and} \quad \mathbb{A}_i = \mathbb{I}_{n_i} - \frac{1}{n_i} \mathbf{1}_{n_i} \mathbf{1}_{n_i}^\otimes.$$

Proposition 6.1. Let \mathbf{X} be a random vector of length n with mean $\mathbb{E} \mathbf{X} = \boldsymbol{\mu}$ and variance matrix $\text{var} \mathbf{X} = \mathbb{V}$. Let \mathbb{A} be any $n \times n$ matrix. Then

$$\mathbb{E} \mathbf{X}^\top \mathbb{A} \mathbf{X} = \boldsymbol{\mu}^\top \mathbb{A} \boldsymbol{\mu} + \text{tr} \mathbb{A} \mathbb{V}.$$

According to Proposition 6.1,

$$\mathbb{E} SS_A = \mathbb{E} \mathbf{Z}^\top \mathbb{A} \mathbf{Z} = (I - 1) \sigma_e^2 + \sum_{i=1}^I n_i (\mu_i - \bar{\mu})^2,$$

where $\bar{\mu} = n^{-1} \sum n_i \mu_i = \mu + n^{-1} \sum n_i \beta_i = \mu + \bar{\beta}$, and

$$\mathbb{E} SS_e = \sum_{i=1}^I \mathbb{E} \mathbf{z}_i^\top \mathbb{A}_i \mathbf{z}_i = (n - I) \sigma_e^2.$$

Hence the expectations of mean squares (Table 6.1, column EMS) are

$$\begin{aligned} \mathbb{E} MS_A &= \sigma_e^2 + \frac{1}{I-1} \sum_{i=1}^I n_i (\mu_i - \bar{\mu})^2 = \sigma_e^2 + \frac{1}{I-1} \sum_{i=1}^I n_i (\beta_i - \bar{\beta})^2 = \sigma_e^2 + Q_\beta, \\ \mathbb{E} MS_e &= \sigma_e^2, \end{aligned}$$

not assuming that H_0 holds and that the data are normal. $\mathbb{E} MS_A$ is the sum of the residual variance σ_e^2 and weighted between-subject variability (sample variance of β_i).

When H_0 is true, $\mathbb{E} MS_A = \mathbb{E} MS_e$ and, under normality, $F = MS_A/MS_e$ has an F distribution with $I - 1$ and $n - I$ degrees of freedom.

6.1.2 One-way ANOVA with random effects

Suppose that the subjects (levels of factor A) represent a random sample from some general population of subjects. Then we may not want to test differences in the means of the I particular subjects participating in the experiment (as classical one-way ANOVA does) but to test whether $\mathbb{E} Y_{ij}$ depends on the subject in the context of all subjects belonging to the whole general population. To do this, we need to acknowledge in the model that we may have taken different subjects who would have had different means than the I subjects we were observing. A natural way to do this is to consider the subject effects random rather than fixed.

Assume

$$Y_{ij} = \mu + b_i + \varepsilon_{ij},$$

where μ is the overall mean, ε_{ij} are iid error terms satisfying $\mathbb{E} \varepsilon_{ij} = 0$, $\text{var} \varepsilon_{ij} = \sigma_e^2$, and b_1, \dots, b_I are iid random variables with $\mathbb{E} b_i = 0$, $\text{var} b_i = \sigma_b^2$. We also assume independence between (b_1, \dots, b_I) and $(\varepsilon_{11}, \dots, \varepsilon_{In_I})$. The random variables b_i are called *random effects**.

There are 3 unknown parameters in this model, in particular $\mu, \sigma_b^2, \sigma_e^2$. The parameter σ_b^2 expresses how different are the subjects in the general populations in terms of the subject-specific mean response.

It is easy to calculate the conditional (given subject) and unconditional (in the general population) moments of Y_{ij} :

$$\begin{aligned} \mathbb{E}[Y_{ij} | b_i] &= \mu + b_i, & \mathbb{E} Y_{ij} &= \mu, \\ \text{var}[Y_{ij} | b_i] &= \sigma_e^2, & \text{var} Y_{ij} &= \sigma_b^2 + \sigma_e^2. \end{aligned}$$

The parameters σ_b^2 and σ_e^2 are called *variance components*† because the total variability in Y_{ij} is decomposed into variability between subjects σ_b^2 (how different are the subjects from each other) and variability within subjects σ_e^2 (how different are the observations made on the same subject from each other).

* Český náhodné efekty † Český komponenty rozptylu

Unlike in classical one-way ANOVA, the observations are not independent:

$$\begin{aligned} \text{cov}(Y_{ij}, Y_{i'j'}) &= 0 && \text{for } i \neq i', \\ \text{cov}(Y_{ij}, Y_{ij'}) &= \text{var } b_i = \sigma_b^2 && \text{for } j \neq j', \\ \text{cor}(Y_{ij}, Y_{ij'}) &= \frac{\sigma_b^2}{\sigma_b^2 + \sigma_e^2} \geq 0 && \text{for } j \neq j'. \end{aligned}$$

The observations are correlated within each subject unless $\sigma_b^2 = 0$. The data has group-dependent structure. The correlations arise because the differences between subjects are not expressed in the mean structure (μ) but added to the random component of the model ($b_i + \varepsilon_{ij}$). In classical ANOVA, the mean structure ($\mu + \beta_i$) captures the differences between subjects and the random component of the model (ε_{ij}) consists of independent variables.*

We would like to test the hypothesis that all subjects in the general population have the same mean response. Within our model, this hypothesis can be expressed as

$$H_0 : \sigma_b^2 = 0 \quad \text{against} \quad H_1 : \sigma_b^2 > 0.$$

The analysis proceeds exactly as in the classical one-way ANOVA. Define SS_A , SS_e , MS_A and MS_e as in the previous section. It can be shown using Proposition 6.1 that

$$\begin{aligned} \mathbb{E} MS_A &= \sigma_e^2 + \frac{1}{I-1} \left(n - \frac{\sum n_i^2}{n} \right) \sigma_b^2, \\ \mathbb{E} MS_e &= \sigma_e^2. \end{aligned}$$

If the null hypothesis is true then $\mathbb{E} MS_A = \mathbb{E} MS_e = \sigma_e^2$. If the null hypothesis is true and normality holds for both ε_{ij} and b_i then $F = MS_A/MS_e$ has an F distribution with $I - 1$ and $n - I$ degrees of freedom. Thus, the analysis proceeds according to Table 6.1 but the null hypothesis and the interpretation of the results are different.

The expected mean squares allow us to find an unbiased estimator of σ_b^2 quite easily. In particular,

$$\hat{\sigma}_b^2 = (MS_A - MS_e) \frac{I-1}{n - \frac{\sum n_i^2}{n}}$$

is unbiased. However, the probability of $MS_A < MS_e$ is positive so this estimator can sometimes attain negative values. We can take $\max(0, \hat{\sigma}_b^2)$ as an estimator but it is no longer unbiased.

* This an illustration of the fact that difference between dependence and independence can be a matter of point of view. The same observations can be considered independent by one analyst and dependent by another analyst without any of them being wrong.

6.1.3 Two-way ANOVA with random effects

Now consider two classification factors A and B, both with randomly selected levels from some general populations. Let

$$Y_{ijk} = \mu + b_i^A + b_j^B + b_{ij}^{AB} + \varepsilon_{ijk},$$

where $i = 1, \dots, I$ are levels of factor A, $j = 1, \dots, J$ are levels of factor B, and $k = 1, \dots, m$ indexes observations made for each combination of the levels of A and B. The design is balanced (equal number of observations for each combination of i and j). The total number of observations is $n = mIJ$.

The parameter μ is the overall mean (a constant), b_i^A , b_j^B , and b_{ij}^{AB} are mutually independent random effects* satisfying $\mathbf{E} b_i^A = 0$, $\text{var} b_i^A = \sigma_A^2$, $\mathbf{E} b_j^B = 0$, $\text{var} b_j^B = \sigma_B^2$, and $\mathbf{E} b_{ij}^{AB} = 0$, $\text{var} b_{ij}^{AB} = \sigma_{AB}^2$. Finally, ε_{ijk} are iid error terms satisfying $\mathbf{E} \varepsilon_{ijk} = 0$, $\text{var} \varepsilon_{ijk} = \sigma_e^2$, independent from all the random effects.

There are 5 unknown parameters in this model, in particular $\mu, \sigma_A^2, \sigma_B^2, \sigma_{AB}^2, \sigma_e^2$.

We get

$$\begin{aligned} \mathbf{E}[Y_{ijk} | b_i^A, b_j^B, b_{ij}^{AB}] &= \mu + b_i^A + b_j^B + b_{ij}^{AB} \\ \text{var}[Y_{ijk} | b_i^A, b_j^B, b_{ij}^{AB}] &= \sigma_e^2 \\ \mathbf{E} Y_{ijk} &= \mu \\ \text{var} Y_{ijk} &= \sigma_A^2 + \sigma_B^2 + \sigma_{AB}^2 + \sigma_e^2. \end{aligned}$$

The total variability in Y_{ijk} is decomposed into four variance components.

The covariances depend on how many levels of A and B are shared by the observations:

$$\begin{aligned} \text{cov}(Y_{ijk}, Y_{ijl}) &= \sigma_A^2 + \sigma_B^2 + \sigma_{AB}^2 && \text{for } k \neq l, \\ \text{cov}(Y_{ijk}, Y_{ij'l}) &= \sigma_A^2 < \text{cov}(Y_{ijk}, Y_{ijl}) && \text{for } j \neq j', \\ \text{cov}(Y_{ijk}, Y_{i'jl}) &= \sigma_B^2 < \text{cov}(Y_{ijk}, Y_{ijl}) && \text{for } i \neq i', \\ \text{cov}(Y_{ijk}, Y_{i'j'l}) &= 0 && \text{for } i \neq i', j \neq j'. \end{aligned}$$

The covariance (and correlation) is largest between two observations made on the same levels of both A and B. It is smaller between two observations made on the same level of either A or B (not both). Two observations made on different levels of both A and B are independent. This data does not have group-dependent structure, there are no mutually independent groups.

* The random effects are independent from each other as well as from random effects pertaining to other levels of A and B

Define the sums of squares and mean squares exactly as in the balanced two-way ANOVA problem with fixed effects:

$$\begin{aligned}
 SS_A &= mJ \sum_{i=1}^I (\bar{Y}_{i..} - \bar{Y}_{...})^2, & MS_A &= \frac{SS_A}{I-1}, \\
 SS_B &= mI \sum_{j=1}^J (\bar{Y}_{.j.} - \bar{Y}_{...})^2, & MS_B &= \frac{SS_B}{J-1}, \\
 SS_{AB} &= m \sum_{i=1}^I \sum_{j=1}^J (\bar{Y}_{ij.} - \bar{Y}_{i..} - \bar{Y}_{.j.} + \bar{Y}_{...})^2, & MS_{AB} &= \frac{SS_{AB}}{(I-1)(J-1)}, \\
 SS_e &= \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^m (Y_{ijk} - \bar{Y}_{ij.})^2, & MS_e &= \frac{SS_e}{(m-1)IJ}.
 \end{aligned}$$

Rewriting the sums of squares as quadratic forms and using Proposition 6.1, we get expected mean squares as follows:

$$\begin{aligned}
 \mathbf{E} MS_A &= mJ\sigma_A^2 + m\sigma_{AB}^2 + \sigma_e^2, \\
 \mathbf{E} MS_B &= mI\sigma_B^2 + m\sigma_{AB}^2 + \sigma_e^2, \\
 \mathbf{E} MS_{AB} &= m\sigma_{AB}^2 + \sigma_e^2, \\
 \mathbf{E} MS_e &= \sigma_e^2.
 \end{aligned} \tag{6.1}$$

Consider testing the null hypothesis $H_0 : \sigma_A^2 = 0$. The only mean squares term that involves the tested parameter is MS_A . Under the null hypothesis, $\mathbf{E} MS_A = \mathbf{E} MS_{AB}$ and the test statistic is $F_A = MS_A/MS_{AB}$, distributed as $F_{I-1, (I-1)(J-1)}$ under normality and validity of the null hypothesis. This is different from two-way ANOVA with fixed effects, where the test statistic for testing factor A main effects is MS_A/MS_e , distributed as $F_{I-1, (m-1)IJ}$ under normality and validity of the null hypothesis. So here, the test for random effects is different from the test for fixed effects. Because $\mathbf{E} MS_{AB} \geq \mathbf{E} MS_e$, the test usually results in smaller test statistic than the test for fixed effects and rejects less frequently. Thus, in more complex models than simple one-way ANOVA, it is more difficult to demonstrate violation of the null hypothesis when the subjects represent a random sample than if they are considered fixed.

The null hypothesis $H_0 : \sigma_B^2 = 0$ is tested by the statistic $F_B = MS_B/MS_{AB}$, distributed as $F_{J-1, (I-1)(J-1)}$ under normality and validity of the null hypothesis. The null hypothesis $H_0 : \sigma_{AB}^2 = 0$ is tested by the statistic $F_{AB} = MS_{AB}/MS_e$, distributed as $F_{(I-1)(J-1), (m-1)IJ}$ under normality and validity of the null hypothesis. This is the only test in this model that is the same with random effects as with fixed effects.

The system of equations (6.1) is linear in the parameters. It can be used to derive

these unbiased moment estimators of the variance components:

$$\begin{aligned}\hat{\sigma}_e^2 &= MS_e, \\ \hat{\sigma}_{AB}^2 &= \frac{1}{m}(MS_{AB} - MS_e), \\ \hat{\sigma}_A^2 &= \frac{1}{mJ}(MS_A - MS_{AB}), \\ \hat{\sigma}_B^2 &= \frac{1}{mI}(MS_B - MS_{AB}).\end{aligned}$$

All these estimators except $\hat{\sigma}_e^2$ may be negative.

The idea of generalizing the results known from classical ANOVA theory to random effects can be followed for more complex models than this one. However, the results quickly become too difficult to derive and apply. Even the two-way ANOVA model with unbalanced data (number of observations n_{ij} vary with i and j) results in very complicated expressions for EMS . The F-test statistics get more complex and do not have exact F distribution even under normality. Various approximations to the denominator degrees of freedom are used to find approximate F distributions of these statistics (one of them is called *Satterthwaite's approximation*).

6.1.4 Random intercept and slope

Let Y_{ij} be measured on a subject i at the time t_{ij} , $i = 1, \dots, K$, $j = 1, \dots, n_i$. The number of observations on a subject and the timing of the observations may vary. Suppose that EY_{ij} depends linearly on time.

Consider the model

$$Y_{ij} = \beta_0 + \beta_1 t_{ij} + \varepsilon_{ij}, \tag{6.2}$$

where ε_{ij} are independent, $E\varepsilon_{ij} = 0$, $\text{var}\varepsilon_{ij} = \sigma_e^2$. This is the classical linear model. It assures that EY_{ij} is a linear function of time but implies that Y_{i1}, \dots, Y_{in_i} are mutually independent. This is not a realistic assumption for observations that are measured sequentially on the same subject.

Now let us extend the model so that different subjects follow different lines. Let b_i^0 be the deviation of the intercept of the i -th subject from the population intercept β_0 and let b_i^1 be the deviation of the slope of the i -th subject from the population slope β_1 . We could include these parameters as fixed effects into model (6.2). However, if the number of subjects grew to infinity the number of subject-specific parameters would grow also and we would not be able to find their consistent estimators. It is better to view the subjects as a random sample coming from some general population of subjects and to consider their parameters random.

This consideration leads us to the model

$$Y_{ij} = \beta_0 + b_i^0 + (\beta_1 + b_i^1)t_{ij} + \varepsilon_{ij},$$

where $\mathbf{E}\varepsilon_{ij} = 0$, $\text{var}\varepsilon_{ij} = \sigma_e^2$, $\mathbf{E}b_i^0 = 0$, $\text{var}b_i^0 = \sigma_0^2$, $\mathbf{E}b_i^1 = 0$, $\text{var}b_i^1 = \sigma_1^2$, ε_{ij} are independent, (b_i^0, b_i^1) are independent, ε_{ij} and (b_i^0, b_i^1) are independent, and $\text{cov}(b_i^0, b_i^1) = \sigma_{01}$.

Separating the fixed and random part, we get

$$Y_{ij} = \underbrace{\beta_0 + \beta_1 t_{ij}}_{\text{fixed}} + \underbrace{b_i^0 + b_i^1 t_{ij}}_{\text{random}} + \varepsilon_{ij}. \quad (6.3)$$

The fixed part of (6.3) is the same as in the linear model (6.2). However, the random part is richer, it involves time and induces correlations between measurements taken on the same subject.

We have $\mathbf{E}[Y_{ij} | b_i^0, b_i^1] = \beta_0 + b_i^0 + (\beta_1 + b_i^1)t_{ij}$, hence $\mathbf{E}Y_{ij} = \beta_0 + \beta_1 t_{ij}$. The linear dependence on time is preserved, with the same parameters as in model (6.2). Also, $\text{var}[Y_{ij} | b_i^0, b_i^1] = \sigma_e^2$ and hence

$$\text{var}Y_{ij} = \sigma_e^2 + \text{var}(b_i^0 + b_i^1 t_{ij}) = \sigma_e^2 + \sigma_0^2 + 2t_{ij}\sigma_{01} + t_{ij}^2\sigma_1^2.$$

So, under model (6.3), the population variance $\text{var}Y_{ij}$ is a quadratic function of time. Also,

$$\text{cov}(Y_{ij}, Y_{ik}) = \text{cov}(b_i^0 + b_i^1 t_{ij}, b_i^0 + b_i^1 t_{ik}) = \sigma_0^2 + (t_{ij} + t_{ik})\sigma_{01} + t_{ij}t_{ik}\sigma_1^2.$$

It can be shown that $\text{cov}(Y_{ij}, Y_{ik}) \geq (\sigma_0 - t_{ij}\sigma_1)(\sigma_0 - t_{ik}\sigma_1)$.

This model can be written as

$$\mathbf{Y}_i = \mathbb{X}_i\boldsymbol{\beta} + \mathbb{Z}_i\mathbf{b}_i + \boldsymbol{\varepsilon}_i, \quad i = 1, \dots, K$$

where $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^\top$, $\boldsymbol{\varepsilon}_i = (\varepsilon_{i1}, \dots, \varepsilon_{in_i})^\top$,

$$\mathbb{X}_i = \mathbb{Z}_i = \begin{pmatrix} 1 & t_{i1} \\ \vdots & \vdots \\ 1 & t_{in_i} \end{pmatrix},$$

$\boldsymbol{\beta} = (\beta_0, \beta_1)^\top$, $\mathbf{b}_i = (b_i^0, b_i^1)^\top$, $\mathbf{E}\mathbf{b}_i = \mathbf{0}$, $\text{var}\mathbf{b}_i = \boldsymbol{\Psi} = \begin{pmatrix} \sigma_0^2 & \sigma_{01} \\ \sigma_{01} & \sigma_1^2 \end{pmatrix}$, $\mathbf{E}\boldsymbol{\varepsilon}_i = \mathbf{0}$, and $\text{var}\boldsymbol{\varepsilon}_i = \sigma_e^2\mathbb{I}_{n_i}$.

This is an example of a linear mixed effects model to be studied in the next sections.

6.2 Definition of Linear Mixed Effects Model

6.2.1 Single-level LME model

Consider independent subjects $i = 1, \dots, K$ and denote $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^\top$ and $\boldsymbol{\varepsilon}_i = (\varepsilon_{i1}, \dots, \varepsilon_{in_i})^\top$. Let \mathbb{X}_i be an $n_i \times p$ regression matrix for fixed effects, $\boldsymbol{\beta}$ a p -vector of fixed effects, let \mathbb{Z}_i be an $n_i \times q$ regression matrix for random effects.

Definition 6.1. The responses $\mathbf{Y}_1, \dots, \mathbf{Y}_K$ satisfy a (single-level) *linear mixed effects model** [LME model] if

$$\mathbf{Y}_i = \mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i + \boldsymbol{\varepsilon}_i, \quad i = 1, \dots, K, \quad (6.4)$$

where \mathbf{b}_i (*the random effects†*) are independent vectors,

$$\mathbf{b}_i \sim \mathbf{N}_q(\mathbf{0}, \mathbb{D}),$$

$\boldsymbol{\varepsilon}_i$ are independent vectors,

$$\boldsymbol{\varepsilon}_i \sim \mathbf{N}_{n_i}(\mathbf{0}, \sigma_e^2 \mathbb{I}_{n_i}),$$

and $(\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_K)^\top$ is independent of $(\mathbf{b}_1, \dots, \mathbf{b}_K)^\top$.

This model was proposed by Laird and Ware (1982). Its mean structure is described by the term $\mathbb{X}_i \boldsymbol{\beta}$, the random structure is given by $\mathbb{Z}_i \mathbf{b}_i + \boldsymbol{\varepsilon}_i$. Obviously, the random part has expectation $\mathbf{0}$ and variance $\mathbb{Z}_i \mathbb{D} \mathbb{Z}_i^\top + \sigma_e^2 \mathbb{I}_{n_i}$. The assumption $\text{var } \boldsymbol{\varepsilon}_i = \sigma_e^2 \mathbb{I}_{n_i}$ will be relaxed in Section 6.5.

The parameters of the LME model are the fixed effects $\boldsymbol{\beta}$, the variance of the residual term σ_e^2 , and the covariance matrix of random effects \mathbb{D} , a symmetric positive definite matrix. For computational reasons, it is convenient to replace \mathbb{D} by another matrix Δ such that $\Delta^\top \Delta = \sigma_e^2 \mathbb{D}^{-1}$. The matrix Δ (which is not unique) is called *the relative precision factor*. It can be obtained, for example, by Choleski decomposition of a positive definite matrix‡.

Model (6.4) can be also written as

$$\mathbf{Y}_i \sim \mathbf{N}_{n_i}(\mathbb{X}_i \boldsymbol{\beta}, \sigma_e^2 \Sigma_i), \quad (6.5)$$

where $\Sigma_i = \mathbb{Z}_i \mathbb{D} \mathbb{Z}_i^\top / \sigma_e^2 + \mathbb{I}_{n_i}$. This is the marginal form of the LME model, which hides the random effects structure. The marginal model (6.5) cannot be always written as a mixed effects model (6.4) – for example, if the matrix \mathbb{D} is not positive definite. The marginal form can be extended to include all the observations

$$\mathbf{Y} \sim \mathbf{N}_n(\mathbb{X} \boldsymbol{\beta}, \sigma_e^2 \Sigma),$$

where $n = \sum_{i=1}^K n_i$,

$$\mathbf{Y} = \begin{pmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_K \end{pmatrix}, \quad \mathbb{X} = \begin{pmatrix} \mathbb{X}_1 \\ \vdots \\ \mathbb{X}_K \end{pmatrix},$$

and Σ is a block-diagonal matrix with diagonal blocks $\Sigma_1, \dots, \Sigma_K$.

* Český lineární smíšený model † Český náhodné efekty ‡ $\mathbb{A} > 0 \Rightarrow$ there exists a lower triangular \mathbb{L} s.t. $\mathbb{A} = \mathbb{L}\mathbb{L}^\top$.

6.2.2 Multi-level LME model

Another grouping level can be added to the LME as follows. Consider primary groups $i = 1, \dots, K$, and secondary groups $j = 1, \dots, m_i$ nested within the primary groups (such as family j within town i). Observe a multivariate vector $\mathbf{Y}_{ij} = (Y_{ij1}, \dots, Y_{ijn_{ij}})^\top$ for each combination of grouping levels i and j .

Definition 6.2. The independent vectors $\mathbf{Y}_{11}, \dots, \mathbf{Y}_{Km_K}$ satisfy a two-level linear mixed effects model if

$$\mathbf{Y}_{ij} = \mathbb{X}_{ij}\boldsymbol{\beta} + \mathbb{Z}_{i,j}\mathbf{b}_i + \mathbb{Z}_{ij}\mathbf{b}_{ij} + \boldsymbol{\varepsilon}_{ij}, \quad i = 1, \dots, K, \quad j = 1, \dots, m_i, \quad (6.6)$$

where \mathbf{b}_i are independent first-level random effects

$$\mathbf{b}_i \sim \mathbf{N}_{q_1}(\mathbf{0}, \mathbb{D}_1),$$

\mathbf{b}_{ij} are independent second-level random effects

$$\mathbf{b}_{ij} \sim \mathbf{N}_{q_2}(\mathbf{0}, \mathbb{D}_2),$$

$\boldsymbol{\varepsilon}_{ij}$ are independent residual terms

$$\boldsymbol{\varepsilon}_{ij} \sim \mathbf{N}_{n_{ij}}(\mathbf{0}, \sigma_e^2 \mathbb{I}_{n_{ij}})$$

(all the random effects and residual terms are mutually independent), \mathbb{X}_{ij} is an $n_{ij} \times p$ regression matrix for fixed effects, $\boldsymbol{\beta}$ is a p -vector of fixed effects, $\mathbb{Z}_{i,j}$ is an $n_{ij} \times q_1$ regression matrix for first-level random effects, and \mathbb{Z}_{ij} is an $n_{ij} \times q_2$ regression matrix for second-level random effects.

This model can be further extended to three or more levels (*multilevel modelling*). In the next sections, however, we will only present results for single-level LME. Extensions to multi-level models are possible but the notation gets a bit more complex.

6.3 Parameter Estimation

Consider a single-level LME model satisfying Definition 6.1. Suppose there exists a matrix Δ such that $\Delta^\top \Delta = \sigma_e^2 \mathbb{D}^{-1}$. Let $\boldsymbol{\theta}$ be an unconstrained parameter that uniquely determines the matrix Δ (and, together with σ_e^2 , also \mathbb{D}).

6.3.1 Marginal likelihood

Consider the marginal form of the model

$$\mathbf{Y}_i \sim \mathbf{N}_{n_i}(\mathbb{X}_i \boldsymbol{\beta}, \sigma_e^2 \Sigma_i),$$

where

$$\Sigma_i = \mathbb{Z}_i \mathbb{D} \mathbb{Z}_i^\top / \sigma_e^2 + \mathbb{I}_{n_i} = \mathbb{Z}_i (\Delta^\top \Delta)^{-1} \mathbb{Z}_i^\top + \mathbb{I}_{n_i}$$

is a function of parameters $\boldsymbol{\theta}$ only. The log-likelihood can be written as

$$\ell_n(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2) = c - \frac{n}{2} \log \sigma_e^2 - \frac{1}{2} \sum_{i=1}^K \log |\Sigma_i| - \frac{1}{2\sigma_e^2} \sum_{i=1}^K (\mathbf{Y}_i - \mathbb{X}_i \boldsymbol{\beta})^\top \Sigma_i^{-1} (\mathbf{Y}_i - \mathbb{X}_i \boldsymbol{\beta}). \quad (6.7)$$

When $\boldsymbol{\theta}$ and σ_e^2 are fixed, maximizing this with respect to $\boldsymbol{\beta}$ is equivalent to minimizing the weighted least squares criterion

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^K (\mathbf{Y}_i - \mathbb{X}_i \boldsymbol{\beta})^\top \Sigma_i^{-1} (\mathbf{Y}_i - \mathbb{X}_i \boldsymbol{\beta}).$$

The solution is the weighted least squares estimator

$$\hat{\boldsymbol{\beta}}(\boldsymbol{\theta}) = \left(\sum_{i=1}^K \mathbb{X}_i^\top \Sigma_i^{-1} \mathbb{X}_i \right)^{-1} \sum_{i=1}^K \mathbb{X}_i^\top \Sigma_i^{-1} \mathbf{Y}_i.$$

This estimator depends on $\boldsymbol{\theta}$ but not on σ_e^2 . Next, maximize the profile likelihood $\ell_n(\hat{\boldsymbol{\beta}}(\boldsymbol{\theta}), \boldsymbol{\theta}, \sigma_e^2)$ with respect to σ_e^2 given $\boldsymbol{\theta}$. The estimator is obtained by the normalized residual weighted least squares

$$\hat{\sigma}_e^2(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^K (\mathbf{Y}_i - \mathbb{X}_i \hat{\boldsymbol{\beta}}(\boldsymbol{\theta}))^\top \Sigma_i^{-1} (\mathbf{Y}_i - \mathbb{X}_i \hat{\boldsymbol{\beta}}(\boldsymbol{\theta})).$$

Finally, take $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$ and $\hat{\sigma}_e^2(\boldsymbol{\theta})$, plug them into ℓ_n to get profile likelihood for $\boldsymbol{\theta}$

$$\ell_n^*(\boldsymbol{\theta}) = \ell_n(\hat{\boldsymbol{\beta}}(\boldsymbol{\theta}), \boldsymbol{\theta}, \hat{\sigma}_e^2(\boldsymbol{\theta})),$$

and maximize it over $\boldsymbol{\theta}$. The form of this profile likelihood is very complicated; analytical calculation of the score statistic and information matrix is possible only in simple special cases. The problem can be handled by numerical optimization methods that do not rely on analytical derivatives.

In the next sections, we will develop a different approach to parameter estimation in LME models that takes advantage of the random effects structure combined with a clever decomposition of the log-likelihood.

6.3.2 Henderson's mixed model equations

Write the model in the more general form

$$\mathbf{Y} = \mathbb{X}\boldsymbol{\beta} + \mathbb{Z}\mathbf{b} + \boldsymbol{\varepsilon},$$

where \mathbb{X} has n rows and p columns, \mathbb{Z} has n rows and q^* columns, $\mathbf{b} \sim \mathbf{N}_{q^*}(\mathbf{0}, \mathbb{D}_*)$, $\boldsymbol{\varepsilon} \sim \mathbf{N}_n(\mathbf{0}, \Lambda_*)$, $\text{cov}(\mathbf{b}, \boldsymbol{\varepsilon}) = 0$, and $\Sigma = \text{var } \mathbf{Y} = \mathbb{Z}\mathbb{D}_*\mathbb{Z}^\top + \Lambda_*$.

The single-level LME model can be obtained as a special case with

$$q^* = Kq, \quad n = \sum_{i=1}^K n_i, \quad \mathbf{Y} = \begin{pmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_K \end{pmatrix}, \quad \mathbb{X} = \begin{pmatrix} \mathbb{X}_1 \\ \vdots \\ \mathbb{X}_K \end{pmatrix}, \quad \mathbb{Z} = \begin{pmatrix} \mathbb{Z}_1 & 0 & \cdots & 0 \\ 0 & \mathbb{Z}_2 & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbb{Z}_K \end{pmatrix},$$

$$\mathbf{b} = \begin{pmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_K \end{pmatrix}, \quad \mathbb{D}_* = \begin{pmatrix} \mathbb{D} & 0 & \cdots & 0 \\ 0 & \mathbb{D} & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbb{D} \end{pmatrix}, \quad \Lambda_* = \sigma_e^2 \mathbb{I}_n.$$

Consider \mathbb{D}_* and Λ_* known and write the joint density of (\mathbf{Y}, \mathbf{b}) with parameters $\boldsymbol{\beta}$:

$$f(\mathbf{y}, \mathbf{b}; \boldsymbol{\beta}) = f(\mathbf{y} | \mathbf{b}; \boldsymbol{\beta})f(\mathbf{b}) = \frac{1}{(2\pi)^{n/2} |\Lambda_*|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{y} - \mathbb{X}\boldsymbol{\beta} - \mathbb{Z}\mathbf{b})^\top \Lambda_*^{-1}(\mathbf{y} - \mathbb{X}\boldsymbol{\beta} - \mathbb{Z}\mathbf{b})\right\} \\ \times \frac{1}{(2\pi)^{q^*/2} |\mathbb{D}_*|^{1/2}} \exp\left\{-\frac{1}{2}\mathbf{b}^\top \mathbb{D}_*^{-1}\mathbf{b}\right\}.$$

Treat this density as a likelihood for unknown parameters $(\boldsymbol{\beta}, \mathbf{b})$ and maximize it jointly to obtain estimators of both $\boldsymbol{\beta}$ and \mathbf{b} . We get

$$(\hat{\boldsymbol{\beta}}, \hat{\mathbf{b}}) = \arg \min_{\boldsymbol{\beta}, \mathbf{b}} \begin{pmatrix} \mathbf{Y} - \mathbb{X}\boldsymbol{\beta} - \mathbb{Z}\mathbf{b} \\ \mathbf{b} \end{pmatrix}^\top \begin{pmatrix} \Lambda_*^{-1} & 0 \\ 0 & \mathbb{D}_*^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{Y} - \mathbb{X}\boldsymbol{\beta} - \mathbb{Z}\mathbf{b} \\ \mathbf{b} \end{pmatrix}.$$

This can be rewritten as weighted sum of squares

$$\arg \min_{\boldsymbol{\beta}, \mathbf{b}} \left[\begin{pmatrix} \mathbf{Y} \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \mathbb{X} & \mathbb{Z} \\ 0 & -\mathbb{I}_{q^*} \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta} \\ \mathbf{b} \end{pmatrix} \right]^\top \begin{pmatrix} \Lambda_*^{-1} & 0 \\ 0 & \mathbb{D}_*^{-1} \end{pmatrix} \left[\begin{pmatrix} \mathbf{Y} \\ \mathbf{0} \end{pmatrix} - \begin{pmatrix} \mathbb{X} & \mathbb{Z} \\ 0 & -\mathbb{I}_{q^*} \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta} \\ \mathbf{b} \end{pmatrix} \right].$$

This is a weighted least squares problem. The estimators satisfy the system of equations

$$\begin{pmatrix} \mathbb{X}^\top \Lambda_*^{-1} \mathbb{X} & \mathbb{X}^\top \Lambda_*^{-1} \mathbb{Z} \\ \mathbb{Z}^\top \Lambda_*^{-1} \mathbb{X} & \mathbb{Z}^\top \Lambda_*^{-1} \mathbb{Z} + \mathbb{D}_*^{-1} \end{pmatrix} \begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\mathbf{b}} \end{pmatrix} = \begin{pmatrix} \mathbb{X}^\top \Lambda_*^{-1} \mathbf{Y} \\ \mathbb{Z}^\top \Lambda_*^{-1} \mathbf{Y} \end{pmatrix}. \quad (6.8)$$

Equations (6.8) are called *Henderson's mixed model equations* (Henderson 1984).

If all the inverses exist (which we assume they do), the estimators can be written down explicitly:

$$\hat{\boldsymbol{\beta}} = (\mathbb{X}^\top \Sigma^{-1} \mathbb{X})^{-1} (\mathbb{X}^\top \Sigma^{-1} \mathbf{Y}), \\ \hat{\mathbf{b}} = \mathbb{D}_* \mathbb{Z}^\top \Sigma^{-1} (\mathbf{Y} - \mathbb{X} \hat{\boldsymbol{\beta}}),$$

where $\Sigma = \mathbb{Z}\mathbb{D}_*\mathbb{Z}^\top + \Lambda_*$. In the special case of the single-level LME model, we get

$$\hat{\mathbf{b}}_i = \mathbb{D}\mathbb{Z}_i^\top (\mathbb{Z}_i \mathbb{D}\mathbb{Z}_i^\top + \sigma_e^2 \mathbb{I}_{n_i})^{-1} (\mathbf{Y}_i - \mathbb{X}_i \hat{\boldsymbol{\beta}}).$$

Note.

1. $\hat{\boldsymbol{\beta}}$ is the best linear unbiased estimator (BLUE) and consistent estimator of $\boldsymbol{\beta}$ regardless of the distribution of \mathbf{Y} .
2. $\hat{\mathbf{b}}$ is the best linear unbiased predictor (BLUP) of \mathbf{b} . This means that $\forall \mathbf{c} \in \mathbb{R}^q$ $E(\mathbf{c}^\top \tilde{\mathbf{b}} - \mathbf{c}^\top \mathbf{b})^2$ is minimized by $\hat{\mathbf{b}}$ among all zero-mean linear functions $\tilde{\mathbf{b}}$ of \mathbf{Y} .
3. $\hat{\mathbf{b}}$ is the posterior mean of \mathbf{b} given \mathbf{Y} and $\hat{\boldsymbol{\beta}}$ (the prior mean is $\mathbf{0}$).
4. $\hat{\mathbf{b}}$ is a weighted average of $\mathbf{0}$ (the prior mean) and $\bar{\mathbf{b}}$ estimated by weighted least squares as fixed effects from the same data. $\hat{\mathbf{b}}$ is called *shrinkage estimator* because it shrinks the estimated fixed effects towards zero.

Denote

$$\mathbb{C} \equiv \begin{pmatrix} \mathbb{C}_{11} & \mathbb{C}_{12} \\ \mathbb{C}_{21} & \mathbb{C}_{22} \end{pmatrix} = \begin{pmatrix} \mathbb{X}^\top \Lambda_*^{-1} \mathbb{X} & \mathbb{X}^\top \Lambda_*^{-1} \mathbb{Z} \\ \mathbb{Z}^\top \Lambda_*^{-1} \mathbb{X} & \mathbb{Z}^\top \Lambda_*^{-1} \mathbb{Z} + \mathbb{D}_*^{-1} \end{pmatrix}^{-1}.$$

It can be shown that $\mathbb{C}_{11} = (\mathbb{X}^\top \Sigma^{-1} \mathbb{X})^{-1}$.

Proposition 6.2. (*Henderson 1984*) Suppose \mathbb{D}_* and Λ_* are known. Then

- (i) $\hat{\mathbf{b}}$ maximizes $\text{cor}(\tilde{\mathbf{b}}, \mathbf{b})$ among all linear unbiased predictors $\tilde{\mathbf{b}}$.
- (ii) $E[\mathbf{b} | \hat{\mathbf{b}}] = \hat{\mathbf{b}}$.
- (iii) $\text{var} \mathbb{A} \hat{\boldsymbol{\beta}} = \mathbb{A} \mathbb{C}_{11} \mathbb{A}^\top = \mathbb{A} (\mathbb{X}^\top \Sigma^{-1} \mathbb{X})^{-1} \mathbb{A}^\top$.
- (iv) $\text{cov}(\hat{\boldsymbol{\beta}}, \hat{\mathbf{b}}) = \mathbf{0}$, $\text{cov}(\hat{\boldsymbol{\beta}}, \mathbf{b}) = -\mathbb{C}_{21}$.
- (v) $\text{var} \hat{\mathbf{b}} = \text{cov}(\hat{\mathbf{b}}, \mathbf{b}) = \mathbb{D}_* - \mathbb{C}_{22}$.
- (vi) $\text{var}(\hat{\mathbf{b}} - \mathbf{b}) = \mathbb{C}_{22}$.

Corollary. For a single-level LME model with known \mathbb{D} and σ_e^2 , the following holds:

- (i) $\text{var} \hat{\boldsymbol{\beta}} = (\sum_{i=1}^K \mathbb{X}_i^\top \Sigma_i^{-1} \mathbb{X}_i)^{-1}$,
- (ii) $\text{var} \hat{\mathbf{b}}_i = \mathbb{D} \mathbb{Z}_i^\top \left[\Sigma_i^{-1} - \Sigma_i^{-1} \mathbb{X}_i (\sum_{i=1}^K \mathbb{X}_i^\top \Sigma_i^{-1} \mathbb{X}_i)^{-1} \mathbb{X}_i^\top \Sigma_i^{-1} \right] \mathbb{Z}_i \mathbb{D}$,
- (iii) $\text{var}(\hat{\mathbf{b}}_i - \mathbf{b}_i) = \mathbb{D} - \text{var} \hat{\mathbf{b}}_i$.

6.3.3 Maximum likelihood estimation of variance parameters

The likelihood function for a single-level LME model can be written as

$$L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2 | \mathbf{Y}) = \prod_{i=1}^K f(\mathbf{Y}_i | \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2) = \prod_{i=1}^K \int f(\mathbf{Y}_i | \mathbf{b}_i, \boldsymbol{\beta}, \sigma_e^2) f(\mathbf{b}_i | \boldsymbol{\theta}, \sigma_e^2) d\mathbf{b}_i,$$

where

$$f(\mathbf{Y}_i | \mathbf{b}_i, \boldsymbol{\beta}, \sigma_e^2) = \frac{1}{(2\pi\sigma_e^2)^{n_i/2}} \exp\left\{-\frac{1}{2\sigma_e^2} \|\mathbf{Y}_i - \mathbb{X}_i \boldsymbol{\beta} - \mathbb{Z}_i \mathbf{b}_i\|^2\right\}$$

and

$$\begin{aligned} f(\mathbf{b}_i \mid \boldsymbol{\theta}, \sigma_e^2) &= \frac{1}{(2\pi)^{q/2} |\mathbb{D}|^{1/2}} \exp\left\{-\frac{1}{2} \mathbf{b}_i^\top \mathbb{D}^{-1} \mathbf{b}_i\right\} = \\ &= \frac{1}{(2\pi\sigma_e^2)^{q/2} |\Delta|^{-1}} \exp\left\{-\frac{1}{2\sigma_e^2} \|\Delta \mathbf{b}_i\|^2\right\}, \end{aligned}$$

where Δ is a relative precision factor matrix of size $q \times q$ such that $|\Delta| > 0$ and $\Delta^\top \Delta = \sigma_e^2 \mathbb{D}^{-1}$.

The whole likelihood can be rearranged as

$$L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2 \mid \mathbf{Y}) = \prod_{i=1}^K \frac{|\Delta|}{(2\pi\sigma_e^2)^{n_i/2}} \int \frac{1}{(2\pi\sigma_e^2)^{q/2}} \exp\left\{-\frac{1}{2\sigma_e^2} \left\| \tilde{\mathbf{Y}}_i - \tilde{\mathbf{X}}_i \boldsymbol{\beta} - \tilde{\mathbf{Z}}_i \mathbf{b}_i \right\|^2\right\} d\mathbf{b}_i,$$

where

$$\tilde{\mathbf{Y}}_i = \begin{pmatrix} \mathbf{Y}_i \\ \mathbf{0} \end{pmatrix}, \quad \tilde{\mathbf{X}}_i = \begin{pmatrix} \mathbf{X}_i \\ \mathbf{0} \end{pmatrix}, \quad \text{and} \quad \tilde{\mathbf{Z}}_i = \begin{pmatrix} \mathbf{Z}_i \\ \Delta \end{pmatrix}$$

all have $n_i + q$ rows. For fixed $\boldsymbol{\beta}$, the Euclidean norm $\left\| \tilde{\mathbf{Y}}_i - \tilde{\mathbf{X}}_i \boldsymbol{\beta} - \tilde{\mathbf{Z}}_i \mathbf{b}_i \right\|^2$ is minimized over \mathbf{b}_i for

$$\hat{\mathbf{b}}_i = (\tilde{\mathbf{Z}}_i^\top \tilde{\mathbf{Z}}_i)^{-1} \tilde{\mathbf{Z}}_i^\top (\tilde{\mathbf{Y}}_i - \tilde{\mathbf{X}}_i \boldsymbol{\beta})$$

(see Henderson's equations). Thus

$$\left\| \tilde{\mathbf{Y}}_i - \tilde{\mathbf{X}}_i \boldsymbol{\beta} - \tilde{\mathbf{Z}}_i \mathbf{b}_i \right\|^2 = \left\| \tilde{\mathbf{Y}}_i - \tilde{\mathbf{X}}_i \boldsymbol{\beta} - \tilde{\mathbf{Z}}_i \hat{\mathbf{b}}_i \right\|^2 + (\mathbf{b}_i - \hat{\mathbf{b}}_i)^\top \tilde{\mathbf{Z}}_i^\top \tilde{\mathbf{Z}}_i (\mathbf{b}_i - \hat{\mathbf{b}}_i).$$

The first term does not involve \mathbf{b}_i and

$$\int \frac{1}{(2\pi\sigma_e^2)^{q/2}} \exp\left\{-\frac{1}{2\sigma_e^2} (\mathbf{b}_i - \hat{\mathbf{b}}_i)^\top \tilde{\mathbf{Z}}_i^\top \tilde{\mathbf{Z}}_i (\mathbf{b}_i - \hat{\mathbf{b}}_i)\right\} d\mathbf{b}_i = \frac{1}{\sqrt{|\tilde{\mathbf{Z}}_i^\top \tilde{\mathbf{Z}}_i|}}. \quad (6.9)$$

Because $|\tilde{\mathbf{Z}}_i^\top \tilde{\mathbf{Z}}_i| = |\mathbf{Z}_i^\top \mathbf{Z}_i + \Delta^\top \Delta|$, we get

$$L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2 \mid \mathbf{Y}) = \frac{1}{(2\pi\sigma_e^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma_e^2} \sum_{i=1}^K \left\| \tilde{\mathbf{Y}}_i - \tilde{\mathbf{X}}_i \boldsymbol{\beta} - \tilde{\mathbf{Z}}_i \hat{\mathbf{b}}_i \right\|^2\right\} \prod_{i=1}^K \frac{|\Delta|}{\sqrt{|\mathbf{Z}_i^\top \mathbf{Z}_i + \Delta^\top \Delta|}}.$$

This can be maximized in three steps:

1. For fixed $\boldsymbol{\theta}$ and σ_e^2 , minimize $\sum_{i=1}^K \left\| \tilde{\mathbf{Y}}_i - \tilde{\mathbf{X}}_i \boldsymbol{\beta} - \tilde{\mathbf{Z}}_i \mathbf{b}_i \right\|^2$ jointly with respect to $\boldsymbol{\beta}$ and \mathbf{b}_i . Henderson's equations provide explicit formulae for $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta}, \sigma_e^2)$ and $\hat{\mathbf{b}}_i(\boldsymbol{\theta}, \sigma_e^2)$.

2. Maximize $L(\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}, \sigma_e^2), \boldsymbol{\theta}, \sigma_e^2)$ with respect to σ_e^2 for fixed $\boldsymbol{\theta}$. This step yields

$$\widehat{\sigma}_e^2(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^K \left\| \widetilde{\mathbf{Y}}_i - \widetilde{\mathbf{X}}_i \widehat{\boldsymbol{\beta}} - \widetilde{\mathbf{Z}}_i \widehat{\mathbf{b}}_i \right\|^2.$$

3. Plug $\widehat{\boldsymbol{\beta}}$, $\widehat{\mathbf{b}}_i$, and $\widehat{\sigma}_e^2$ into $L(\cdot | \mathbf{Y})$ and get the profiled likelihood

$$L^*(\boldsymbol{\theta}) = L(\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}), \boldsymbol{\theta}, \widehat{\sigma}_e^2(\boldsymbol{\theta}) | \mathbf{Y}) = \frac{\exp\{-n/2\}}{(2\pi\widehat{\sigma}_e^2(\boldsymbol{\theta}))^{n/2}} \prod_{i=1}^K \frac{|\Delta|}{\sqrt{|\mathbf{Z}_i^T \mathbf{Z}_i + \Delta^T \Delta|}}.$$

The profile log-likelihood is

$$\ell^*(\boldsymbol{\theta}) = \log L^*(\boldsymbol{\theta}) = -\frac{1}{2} \left(n + n \log 2\pi + n \log \widehat{\sigma}_e^2(\boldsymbol{\theta}) - 2K \log |\Delta| + \sum_{i=1}^K \log |\mathbf{Z}_i^T \mathbf{Z}_i + \Delta^T \Delta| \right).$$

This can be simplified further by QR and/or Choleski decompositions. Once the MLE $\widehat{\boldsymbol{\theta}}$ is found, it is plugged into $\widehat{\sigma}_e^2(\widehat{\boldsymbol{\theta}})$ and then into $\widehat{\boldsymbol{\beta}}(\widehat{\boldsymbol{\theta}}, \widehat{\sigma}_e^2)$ and $\widehat{\mathbf{b}}_i(\widehat{\boldsymbol{\theta}}, \widehat{\sigma}_e^2)$. For details, see [Jennrich and Schluchter \(1986\)](#) and [Lindstrom and Bates \(1988\)](#).

6.3.4 Restricted maximum likelihood (REML) estimation of variance parameters

The main idea of REML is to get rid of the fixed effects before the variance parameters are estimated. There are two different ways to derive REML estimators, which lead to the same result. Both derivations are shown for the general case $\mathbf{Y} \sim \mathbf{N}_n(\mathbb{X}\boldsymbol{\beta}, \sigma_e^2 \boldsymbol{\Sigma})$. The single-level LME model is a special case of this model.

Derivation of REML estimators by zero-mean contrasts

The REML method was first proposed by [Patterson and Thompson \(1971\)](#), who took the following strategy to derive it. Take a matrix \mathbb{A} such that $\mathbf{E} \mathbb{A} \mathbf{Y} = \mathbf{0}$ and use the likelihood of $\mathbb{A} \mathbf{Y}$ instead of the likelihood of \mathbf{Y} .

The matrix $\mathbb{I}_n - \mathbb{X}(\mathbb{X}^T \boldsymbol{\Sigma}^{-1} \mathbb{X})^{-1} \mathbb{X}^T \boldsymbol{\Sigma}^{-1}$ generates the residuals of weighted least squares, is idempotent and has rank $n - p$. Take \mathbb{A} a $(n - p) \times n$ matrix such that

$$\mathbb{A}^T \mathbb{A} = \mathbb{I}_n - \mathbb{X}(\mathbb{X}^T \boldsymbol{\Sigma}^{-1} \mathbb{X})^{-1} \mathbb{X}^T \boldsymbol{\Sigma}^{-1} \quad \text{and} \quad \mathbb{A} \mathbb{A}^T = \mathbb{I}_{n-p}.$$

Since $\mathbb{A}^T \mathbb{A} \mathbb{X} = \mathbf{0}$ then also $\mathbb{A} \mathbb{X} = \mathbf{0}$ and $\mathbf{E} \mathbb{A} \mathbf{Y} = \mathbb{A} \mathbb{X} \boldsymbol{\beta} = \mathbf{0}$. We have

$$\mathbb{A} \mathbf{Y} \equiv \mathbf{Y}_A \sim \mathbf{N}_{n-p}(\mathbf{0}, \sigma_e^2 \mathbb{A} \boldsymbol{\Sigma} \mathbb{A}^T).$$

Take $\mathbb{G} = (\mathbb{X}^\top \Sigma^{-1} \mathbb{X})^{-1} \mathbb{X}^\top \Sigma^{-1}$. It is a $p \times n$ matrix such that $\mathbb{G} \mathbb{X} = \mathbb{I}_p$, $\mathbb{G} \mathbf{Y} = \widehat{\boldsymbol{\beta}}$ (the WLS estimator of $\boldsymbol{\beta}$) and

$$\mathbb{G} \mathbf{Y} \equiv \widehat{\boldsymbol{\beta}} \sim \mathbb{N}_p(\boldsymbol{\beta}, \sigma_e^2 (\mathbb{X}^\top \Sigma^{-1} \mathbb{X})^{-1}).$$

The random vector $(\mathbf{Y}_A, \widehat{\boldsymbol{\beta}})$ is a linear one-to-one transformation of \mathbf{Y} , hence it is multivariate normal. Because $\text{cov}(\mathbb{A} \mathbf{Y}, \mathbb{G} \mathbf{Y}) = \mathbf{0}$, the components $\mathbb{A} \mathbf{Y}$ and $\mathbb{G} \mathbf{Y}$ are independent.

Because of independence, the joint density of $f_{\mathbf{Y}_A, \widehat{\boldsymbol{\beta}}}$ of $(\mathbf{Y}_A, \widehat{\boldsymbol{\beta}})$ is the product of the marginal densities $f_{\mathbf{Y}_A}$ and $f_{\widehat{\boldsymbol{\beta}}}$, so we can express $f_{\mathbf{Y}_A} = f_{\mathbf{Y}_A, \widehat{\boldsymbol{\beta}}} / f_{\widehat{\boldsymbol{\beta}}}$. Instead of this ratio, we take $f_{\mathbf{Y}} / f_{\widehat{\boldsymbol{\beta}}}$ and use it as a likelihood function for estimating σ_e^2 and $\boldsymbol{\theta}$.

We have

$$f_{\mathbf{Y}}(\mathbf{y}) = \frac{1}{(2\pi\sigma_e^2)^{n/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2\sigma_e^2}(\mathbf{y} - \mathbb{X}\boldsymbol{\beta})^\top \Sigma^{-1}(\mathbf{y} - \mathbb{X}\boldsymbol{\beta})\right\}. \quad (6.10)$$

and

$$f_{\widehat{\boldsymbol{\beta}}}(\widehat{\boldsymbol{\beta}}) = \frac{1}{(2\pi\sigma_e^2)^{p/2}} \left| \mathbb{X}^\top \Sigma^{-1} \mathbb{X} \right|^{1/2} \exp\left\{-\frac{1}{2\sigma_e^2}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbb{X}^\top \Sigma^{-1} \mathbb{X}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})\right\}.$$

Realizing that $\mathbb{X}^\top \Sigma^{-1}(\mathbf{Y} - \mathbb{X}\widehat{\boldsymbol{\beta}}) = 0$, we get

$$(\mathbf{Y} - \mathbb{X}\boldsymbol{\beta})^\top \Sigma^{-1}(\mathbf{Y} - \mathbb{X}\boldsymbol{\beta}) = (\mathbf{Y} - \mathbb{X}\widehat{\boldsymbol{\beta}})^\top \Sigma^{-1}(\mathbf{Y} - \mathbb{X}\widehat{\boldsymbol{\beta}}) + (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbb{X}^\top \Sigma^{-1} \mathbb{X}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}). \quad (6.11)$$

The ratio $f_{\mathbf{Y}}(\mathbf{Y}) / f_{\widehat{\boldsymbol{\beta}}}(\widehat{\boldsymbol{\beta}})$, simplified by the equality (6.11), does not depend on $\boldsymbol{\beta}$. We call it the restricted (residual) likelihood. It is not a proper likelihood because the ratio of densities is not in general a density.

Definition 6.3. The restricted (REML) likelihood is

$$L_R(\boldsymbol{\theta}, \sigma_e^2) = \frac{1}{(2\pi\sigma_e^2)^{(n-p)/2}} |\Sigma|^{-1/2} \left| \mathbb{X}^\top \Sigma^{-1} \mathbb{X} \right|^{-1/2} \times \exp\left\{-\frac{1}{2\sigma_e^2}(\mathbf{Y} - \mathbb{X}\widehat{\boldsymbol{\beta}})^\top \Sigma^{-1}(\mathbf{Y} - \mathbb{X}\widehat{\boldsymbol{\beta}})\right\}.$$

Note. The REML log-likelihood is

$$\ell_R(\boldsymbol{\theta}, \sigma_e^2) = C - \frac{n-p}{2} \log \sigma_e^2 - \frac{1}{2} \log |\Sigma| - \frac{1}{2} \log \left| \mathbb{X}^\top \Sigma^{-1} \mathbb{X} \right| - \frac{1}{2\sigma_e^2}(\mathbf{Y} - \mathbb{X}\widehat{\boldsymbol{\beta}})^\top \Sigma^{-1}(\mathbf{Y} - \mathbb{X}\widehat{\boldsymbol{\beta}}), \quad (6.12)$$

where Σ and $\widehat{\boldsymbol{\beta}}$ depend on $\boldsymbol{\theta}$. It differs from the marginal log-likelihood (6.7) as follows:

$$\ell_R(\boldsymbol{\theta}, \sigma_e^2) = \ell_n(\widehat{\boldsymbol{\beta}}(\boldsymbol{\theta}), \boldsymbol{\theta}, \sigma_e^2) + \frac{p}{2} \log \sigma_e^2 - \frac{1}{2} \log \left| \mathbb{X}^\top \Sigma^{-1} \mathbb{X} \right| + C_1.$$

Maximizing (6.12) over σ_e^2 for fixed $\boldsymbol{\theta}$ yields the following REML estimator of σ_e^2

$$\widehat{\sigma}_e^2(\boldsymbol{\theta}) = \frac{1}{n-p}(\mathbf{Y} - \mathbb{X}\widehat{\boldsymbol{\beta}})^\top \Sigma^{-1}(\mathbf{Y} - \mathbb{X}\widehat{\boldsymbol{\beta}}).$$

This differs from the MLE by the subtraction of the number of parameters p from the number of observations n in the denominator. Unlike the MLE, the REML estimator is unbiased when $\boldsymbol{\theta}$ is known and agrees with the usual unbiased estimator of residual variance in linear models.

Repeating the steps leading to the derivation of the usual profile likelihood in Section 6.3.3, we end up with REML profile log-likelihood for estimating $\boldsymbol{\theta}$ in the single-level LME model written in the form

$$\begin{aligned} \ell_R^*(\boldsymbol{\theta}) = \ell_R(\boldsymbol{\theta}, \sigma_e^2(\boldsymbol{\theta})) = & -\frac{1}{2} \left[(n-p) \log \widehat{\sigma}_e^2(\boldsymbol{\theta}) - 2K \log |\Delta| + \right. \\ & \left. + \sum_{i=1}^K \log \left| \mathbf{Z}_i^\top \mathbf{Z}_i + \Delta^\top \Delta \right| + \log \left| \sum_{i=1}^K \mathbb{X}_i^\top \Sigma_i^{-1} \mathbb{X}_i \right| \right] + C_1. \end{aligned}$$

This is maximized to obtain the REML estimator $\widehat{\boldsymbol{\theta}}_R$ of $\boldsymbol{\theta}$.

Next, we calculate $\widehat{\sigma}_e^2(\widehat{\boldsymbol{\theta}}_R)$ and, finally, $\widehat{\boldsymbol{\theta}}_R$ is plugged into the weighted least squares estimator to obtain $\widehat{\boldsymbol{\beta}}$.

Derivation of REML estimators by integration over $\boldsymbol{\beta}$

This is a completely different way to derive REML estimators, which nevertheless leads to the same result. It was proposed by Harville (1974).

We want to get rid of the unknown $\boldsymbol{\beta}$ in the likelihood. We can do this by averaging the density of \mathbf{Y} over all possible values of $\boldsymbol{\beta}$ and take the result as a new likelihood for estimation of σ_e^2 and $\boldsymbol{\theta}$.

So the idea is to take

$$L_R(\boldsymbol{\theta}, \sigma_e^2) = \int L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2) d\boldsymbol{\beta}.$$

This can be also viewed in a Bayesian context. Consider $\boldsymbol{\beta}$ random and assign a non-informative improper prior density $h(\boldsymbol{\beta}) = 1$ on \mathbb{R}^p to the random $\boldsymbol{\beta}$. Calculating the expectation of $L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2)$ over $\boldsymbol{\beta}$ yields $L_R(\boldsymbol{\theta}, \sigma_e^2)$.

We will show that this likelihood is the same as the REML likelihood from Definition 6.3. We need to calculate

$$\int L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2) d\boldsymbol{\beta} = \int f_{\mathbf{Y}}(\mathbf{Y}; \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2) d\boldsymbol{\beta},$$

where $f_{\mathbf{Y}}(\mathbf{y})$ is defined by (6.10). From the decomposition (6.11), we get

$$\int f_{\mathbf{Y}}(\mathbf{Y}; \boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2) d\boldsymbol{\beta} = \frac{1}{(2\pi\sigma_e^2)^{n/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2\sigma_e^2}(\mathbf{Y} - \mathbb{X}\hat{\boldsymbol{\beta}})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{Y} - \mathbb{X}\hat{\boldsymbol{\beta}})\right\} \times \int \exp\left\{-\frac{1}{2\sigma_e^2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbb{X}^\top \boldsymbol{\Sigma}^{-1} \mathbb{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\right\} d\boldsymbol{\beta}. \quad (6.13)$$

Using the same trick as in (6.9), we can see that

$$\int \exp\left\{-\frac{1}{2\sigma_e^2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbb{X}^\top \boldsymbol{\Sigma}^{-1} \mathbb{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\right\} d\boldsymbol{\beta} = (2\pi\sigma_e^2)^{q/2} \left|\mathbb{X}^\top \boldsymbol{\Sigma}^{-1} \mathbb{X}\right|^{-1/2}.$$

Thus, the expression calculated from (6.13) is the same as the restricted likelihood in Definition 6.3. Both approaches to deriving REML estimators lead to the same result.

6.3.5 Comparison of REML versus ML estimators

To understand REML estimation, one has to realize three important facts:

- REML likelihood only estimates the variance parameters, not the fixed effects;
- REML likelihood is not a proper likelihood;
- the difference between REML likelihood and ordinary likelihood is asymptotically negligible.

We can summarize the strengths and weaknesses of REML (relative to MLE) as follows:

Advantages of REML

- The variance estimators agree with those used in the analysis of variance and linear regression.
- Under certain conditions (balanced design), REML estimators of variance parameters are unbiased.

Disadvantages of REML

- Because of the extra term $|\mathbb{X}^\top \boldsymbol{\Sigma}^{-1} \mathbb{X}|^{-1/2}$ in REML likelihood, the estimates of σ_e^2 and $\boldsymbol{\theta}$ depend on the parametrization of the fixed part of the model.

Consider $\mathbb{X} \neq \mathbb{X}^*$ such that $\mathcal{M}(\mathbb{X}) = \mathcal{M}(\mathbb{X}^*)$. If MLE is used, then the likelihood, variance parameter estimates, residuals, and fitted values, are the same for both parametrizations. If REML is used, then the likelihood, variance parameter estimates, residuals, and fitted values differ even though the two models are equivalent.

- REML likelihood (or any quantity derived from the likelihood, e.g. AIC or BIC) cannot be used to compare two models with different fixed effects structures or even parametrizations of fixed effects.

6.4 Hypothesis Testing and Confidence Intervals

6.4.1 Asymptotic approach based on MLE theory

It can be shown that in the single-level LME model the terms that express the difference between the regular likelihood and REML likelihood vanish as the number of groups K grows. Thus, the asymptotic theory for ML and REML estimators is the same, with some exceptions explained below.

Theorem 6.3. *If the model holds and regularity assumptions for the MLE theory are satisfied*

- The estimators $\hat{\beta}$, $\hat{\sigma}_e^2$ and $\hat{\theta}$ are consistent as $K \rightarrow \infty$.
- Let r be the dimension of θ . Then

$$\sqrt{K} \begin{pmatrix} \hat{\beta} - \beta \\ \hat{\sigma}_e^2 - \sigma_e^2 \\ \hat{\theta} - \theta \end{pmatrix} \xrightarrow{D} N_{p+r+1}(\mathbf{0}, I^{-1}) \quad \text{as } K \rightarrow \infty,$$

where I is the Fisher information matrix with the following structure

$$I = \begin{pmatrix} I_{\beta} & 0 & 0 \\ 0 & I_{\sigma_e^2} & I_{\sigma_e^2, \theta} \\ 0 & I_{\theta, \sigma_e^2} & I_{\theta} \end{pmatrix}.$$

- Let M be a model, ℓ_M its log-likelihood, and $\hat{\beta}$, $\hat{\sigma}_e^2$, and $\hat{\theta}$ parameter estimates. Let S be a submodel, ℓ_S its log-likelihood, and $\tilde{\beta}$, $\tilde{\sigma}_e^2$, and $\tilde{\theta}$ parameter estimates. If the submodel holds and if the submodel does not specify parameters to be at the boundary of the parameter space then

$$LR = 2[\ell_M(\hat{\beta}, \hat{\sigma}_e^2, \hat{\theta}) - \ell_S(\tilde{\beta}, \tilde{\sigma}_e^2, \tilde{\theta})] \xrightarrow{D} \chi_m^2,$$

where m is the difference in the number of parameters between the model M and the submodel S .

The meaning of Theorem 6.3 needs to be carefully explained in a series of important notes.

Note.

1. The asymptotic results require that the number of independent groups K grows to infinity. It is not enough to keep K bounded and to increase all n_i .

2. **Estimation of fixed effects:**

- $\sqrt{K}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} N_p(\mathbf{0}, I_{\boldsymbol{\beta}}^{-1})$ as $K \rightarrow \infty$ if the model holds, where $I_{\boldsymbol{\beta}} = E \mathbb{X}_i^T \Sigma_i^{-1} \mathbb{X}_i / \sigma_e^2$.
- $\hat{\boldsymbol{\beta}}$ is consistent even if the random part of the model is incorrectly specified but its asymptotic variance is different from $I_{\boldsymbol{\beta}}^{-1}$.
- If the model holds then the asymptotic variance $I_{\boldsymbol{\beta}}^{-1}$ can be consistently estimated by

$$\left[\frac{1}{K} \sum_{i=1}^K \mathbb{X}_i^T (\hat{\sigma}_e^2 \mathbb{I}_{n_i} + \mathbb{Z}_i D(\hat{\boldsymbol{\theta}}) \mathbb{Z}_i^T)^{-1} \mathbb{X}_i \right]^{-1}.$$

3. **Likelihood ratio tests:**

- Theorem 6.3(iii) also holds for REML log-likelihood as long as the fixed part \mathbb{X} is the same in the model as in the submodel.
- Likelihood ratio tests for reduction in the dimension of $\boldsymbol{\beta}$ (tests of the fixed part of the model) are not recommended because of slow convergence to the limiting χ^2 distribution; these tests are more suitable for testing variance components.
- Theorem 6.3(iii) does not hold when the submodel involves a removal of a variance component. This issue is discussed in more detail in the next section.

6.4.2 Likelihood ratio tests

Likelihood ratio tests are not recommended for testing hypotheses about $\boldsymbol{\beta}$ but they are useful for testing hypotheses about variance components (functions of $\boldsymbol{\theta}$). However, variance component testing frequently violates the regularity assumptions of MLE theory, which makes the asymptotic distribution of the LR test statistic specified by Theorem 6.3(iii) invalid.

Consider a model with a single random intercept b_i added to the linear predictor describing the fixed part of the model. Let the variance of b_i be denoted $\sigma_b^2 \geq 0$. The parameter space for the variance component σ_b^2 is not an open set because it includes the zero value. Indeed, testing the presence of the random intercept in the model is equivalent to testing the hypothesis $H_0 : \sigma_b^2 = 0$ against the alternative $H_1 : \sigma_b^2 > 0$. The null hypothesis specifies a value of the parameter σ_b^2 that lies at the border of the parameter space. Because of this, the LR test statistic comparing log-likelihoods with and without the random intercept does not converge to the usual χ_1^2 distribution.

This problem persists in all likelihood ratio tests about a removal of a variance

component (which reduces the dimension of $\boldsymbol{\theta}$ from q to $q - m$). Such test statistics violate the conditions of Theorem 6.3(iii) and do not have asymptotic χ_m^2 distribution.

We provide several examples of the true limiting distributions of likelihood ratio test statistics in important special cases addressed by [Stram and Lee \(1994; 1995\)](#). Proofs are omitted. Recall that $\mathbb{D} = \text{var } \mathbf{b}_i$ and consider a few important hypotheses about \mathbb{D} .

1. $H_0 : \mathbb{D} = 0$ against $H_1 : \mathbb{D} = d_{11} > 0$

This is a test for the presence of a single variance component, the number of tested parameters is $m = 1$. If the null hypothesis is true then $LR \xrightarrow{D} Z$ where the random variable Z is distributed as an equal mixture of the constant 0 and the χ_1^2 distribution.

If the observed value of the LR test statistic is λ the correct asymptotic P-value for the LR test can be calculated as $0.5P[\chi_1^2 \geq \lambda]$.

2. $H_0 : \mathbb{D} = \begin{pmatrix} \mathbb{D}_{11} & \mathbf{0} \\ \mathbf{0}^\top & 0 \end{pmatrix}$ against $H_1 : \mathbb{D} = \begin{pmatrix} \mathbb{D}_{11} & \mathbf{d}_{12} \\ \mathbf{d}_{12}^\top & d_{22} \end{pmatrix}$,

where the dimension of \mathbb{D}_{11} is $q \times q$ and $d_{22} > 0$ is a scalar. This is a test for adding a single variance component to $q > 0$ existing ones, the number of tested parameters is $m = q + 1$. If the null hypothesis is true then $LR \xrightarrow{D} Z$ where the random variable Z is distributed as an equal mixture of the χ_q^2 and the χ_{q+1}^2 distributions.

The correct asymptotic P-value for the LR test is $0.5P[\chi_q^2 \geq \lambda] + 0.5P[\chi_{q+1}^2 \geq \lambda]$.

3. $H_0 : \mathbb{D} = \begin{pmatrix} \mathbb{D}_{11} & \mathbf{0} \\ \mathbf{0}^\top & \mathbf{0} \end{pmatrix}$ against $H_1 : \mathbb{D} = \begin{pmatrix} \mathbb{D}_{11} & \mathbb{D}_{12} \\ \mathbb{D}_{12}^\top & \mathbb{D}_{22} \end{pmatrix}$,

where the dimension of \mathbb{D}_{11} is $q \times q$ and the dimension of $\mathbb{D}_{22} > 0$ is $k \times k$. This is a test for adding k variance components to $q > 0$ existing ones, the number of tested parameters is $m = k(k + 1)/2 + qk$. If the null hypothesis is true then $LR \xrightarrow{D} Z$ where the random variable Z is distributed as an unequal mixture of χ^2 distributions, all of them with degrees of freedom not larger than $k(k + 1)/2 + qk$. The weights in the limiting mixture are rather difficult to calculate.

In all these special cases, the true limiting distribution of the LR statistic is stochastically smaller than the asymptotic distribution χ_m^2 implied by Theorem 6.3(iii). Therefore, the critical values and p -values based on the standard MLE theory are larger than the true critical values and p -values. Likelihood ratio tests of variance components based on standard MLE theory are *conservative*; they provide tests with true level much smaller than the desired α and have much smaller power to detect deviations from the null hypothesis.

6.4.3 t tests and F tests for fixed effects

Suppose Σ is known and take $\hat{\boldsymbol{\beta}} = (\mathbb{X}^T \Sigma^{-1} \mathbb{X})^{-1} (\mathbb{X}^T \Sigma^{-1} \mathbf{Y})$. Consider the problem of estimation and testing linear combinations $\mathbf{c}^T \boldsymbol{\beta}$ of fixed effects for some vector of constants $\mathbf{c} \neq \mathbf{0}$.

Suppose that σ_e^2 is known, too. Then it is easy to show that for any $\mathbf{c} \neq \mathbf{0}$

$$\frac{\mathbf{c}^T \hat{\boldsymbol{\beta}} - \mathbf{c}^T \boldsymbol{\beta}}{\sqrt{\sigma_e^2 \mathbf{c}^T (\mathbb{X}^T \Sigma^{-1} \mathbb{X})^{-1} \mathbf{c}}} \sim \mathbf{N}(0, 1).$$

This could be used to test $H_0 : \mathbf{c}^T \boldsymbol{\beta} = c_0$ and build exact confidence intervals for $\mathbf{c}^T \boldsymbol{\beta}$ if the true variance parameters were known.

Now let σ_e^2 be unknown and use the REML estimator

$$\hat{\sigma}_e^2 = \frac{1}{n-p} (\mathbf{Y} - \mathbb{X} \hat{\boldsymbol{\beta}})^T \Sigma^{-1} (\mathbf{Y} - \mathbb{X} \hat{\boldsymbol{\beta}}).$$

Then it is equally easy to show that for any $\mathbf{c} \neq \mathbf{0}$

$$\frac{\mathbf{c}^T \hat{\boldsymbol{\beta}} - \mathbf{c}^T \boldsymbol{\beta}}{\sqrt{\hat{\sigma}_e^2 \mathbf{c}^T (\mathbb{X}^T \Sigma^{-1} \mathbb{X})^{-1} \mathbf{c}}} \sim t_{n-p}.$$

This could be used to test $H_0 : \mathbf{c}^T \boldsymbol{\beta} = c_0$ and build exact confidence intervals for $\mathbf{c}^T \boldsymbol{\beta}$ if Σ were known and σ_e^2 unknown.

With both Σ and σ_e^2 unknown, take ML or REML estimators $\hat{\Sigma}$ and $\hat{\sigma}_e^2$ of Σ and σ_e^2 , respectively and redefine $\hat{\boldsymbol{\beta}} = (\mathbb{X}^T \hat{\Sigma}^{-1} \mathbb{X})^{-1} (\mathbb{X}^T \hat{\Sigma}^{-1} \mathbf{Y})$. By Theorem 6.3(ii) and the note below it,

$$\frac{\mathbf{c}^T \hat{\boldsymbol{\beta}} - \mathbf{c}^T \boldsymbol{\beta}}{\sqrt{\hat{\sigma}_e^2 \mathbf{c}^T (\mathbb{X}^T \hat{\Sigma}^{-1} \mathbb{X})^{-1} \mathbf{c}}} \xrightarrow{D} \mathbf{N}(0, 1)$$

for any $\mathbf{c} \neq \mathbf{0}$. However, the asymptotics requires $K \rightarrow \infty$ and the approximation may not work well in practical data analyses. Unfortunately, it is no longer possible to obtain the exact small-sample distribution of the left hand side.

There are several options for better approximations of the small-sample distribution than by standard normal:

- Use t_{n-p} as if Σ were known. This is implemented in the R library `nlme`.
- Use a better approximation $t_{\hat{f}}$, where \hat{f} is calculated from the data.
 - Satterthwaite approximation ([Satterthwaite 1946](#))
 - Kenward-Roger approximation ([Kenward and Roger 1997](#))

These approximations should be more precise than t_{n-p} . Kenward-Roger approximation is the preferred one. Both are implemented in the SAS procedure `proc MIXED`.

Now take a fixed $m \times p$ matrix \mathbb{A} , where $m < p$, $r(\mathbb{A}) = m$. Consider the problem of testing the hypothesis $H_0 : \mathbb{A}\boldsymbol{\beta} = \boldsymbol{\gamma}_0$ against $H_1 : \mathbb{A}\boldsymbol{\beta} \neq \boldsymbol{\gamma}_0$, where $\boldsymbol{\gamma}_0$ is some vector of constants, usually zeros.

If σ_e^2 and Σ are known then

$$\mathbb{A}\hat{\boldsymbol{\beta}} \sim \mathbf{N}_m(\mathbb{A}\boldsymbol{\beta}, \sigma_e^2 \mathbb{A}(\mathbb{X}^\top \Sigma^{-1} \mathbb{X})^{-1} \mathbb{A}^\top)$$

and

$$(\mathbb{A}\hat{\boldsymbol{\beta}} - \mathbb{A}\boldsymbol{\beta})^\top [\sigma_e^2 \mathbb{A}(\mathbb{X}^\top \Sigma^{-1} \mathbb{X})^{-1} \mathbb{A}^\top]^{-1} (\mathbb{A}\hat{\boldsymbol{\beta}} - \mathbb{A}\boldsymbol{\beta}) \sim \chi_m^2.$$

If σ_e^2 is unknown and Σ is known then

$$F = \frac{1}{m} (\mathbb{A}\hat{\boldsymbol{\beta}} - \mathbb{A}\boldsymbol{\beta})^\top [\hat{\sigma}_e^2 \mathbb{A}(\mathbb{X}^\top \Sigma^{-1} \mathbb{X})^{-1} \mathbb{A}^\top]^{-1} (\mathbb{A}\hat{\boldsymbol{\beta}} - \mathbb{A}\boldsymbol{\beta}) \sim F_{m, n-p}.$$

If σ_e^2 and Σ are unknown then the statistic F does not have exact $F_{m, n-p}$ distribution but it can be used as an approximation valid for large K . This is implemented in R library `nlme`. Alternatively, one can use an approximation $F_{m, \hat{f}}$, where \hat{f} is calculated from the data (Satterthwaite approximation). However, \hat{f} is different from and more complicated than \hat{f} used in the t-test.

Satterthwaite approximation for t-test

Let us briefly describe how Satterthwaite approximation for t-test can be obtained.

Write

$$\frac{\mathbf{c}^\top \hat{\boldsymbol{\beta}} - \mathbf{c}^\top \boldsymbol{\beta}}{\sqrt{\hat{\sigma}_e^2 \mathbf{c}^\top (\mathbb{X}^\top \hat{\Sigma}^{-1} \mathbb{X})^{-1} \mathbf{c}}} = \frac{\frac{\mathbf{c}^\top \hat{\boldsymbol{\beta}} - \mathbf{c}^\top \boldsymbol{\beta}}{\sqrt{\sigma_e^2 \mathbf{c}^\top (\mathbb{X}^\top \Sigma^{-1} \mathbb{X})^{-1} \mathbf{c}}}}{\sqrt{\frac{f \hat{\sigma}_e^2 \mathbf{c}^\top (\mathbb{X}^\top \hat{\Sigma}^{-1} \mathbb{X})^{-1} \mathbf{c}}{\sigma_e^2 \mathbf{c}^\top (\mathbb{X}^\top \Sigma^{-1} \mathbb{X})^{-1} \mathbf{c}} \frac{1}{f}}}.$$

The numerator has standard normal distribution. To show that the distribution of this statistic can be approximated by χ_f^2 , we seek constants f and ψ such that

$$\frac{f \hat{\sigma}_e^2 \mathbf{c}^\top (\mathbb{X}^\top \hat{\Sigma}^{-1} \mathbb{X})^{-1} \mathbf{c}}{\psi} \sim \chi_f^2.$$

We use the method of moments to find f and ψ from the two equations

$$\begin{aligned} \frac{f}{\psi} \mathbb{E} \hat{\sigma}_e^2 \mathbf{c}^\top (\mathbb{X}^\top \hat{\Sigma}^{-1} \mathbb{X})^{-1} \mathbf{c} &= f, \\ \frac{f^2}{\psi^2} \text{var} \hat{\sigma}_e^2 \mathbf{c}^\top (\mathbb{X}^\top \hat{\Sigma}^{-1} \mathbb{X})^{-1} \mathbf{c} &= 2f. \end{aligned}$$

We get

$$f = \frac{2[\mathbb{E} \hat{\sigma}_e^2 \mathbf{c}^\top (\mathbb{X}^\top \hat{\Sigma}^{-1} \mathbb{X})^{-1} \mathbf{c}]^2}{\text{var} \hat{\sigma}_e^2 \mathbf{c}^\top (\mathbb{X}^\top \hat{\Sigma}^{-1} \mathbb{X})^{-1} \mathbf{c}}.$$

Hence the number of degrees of freedom f is estimated by $\hat{f} = (2\hat{E}^2)/\hat{V}$, where

$$\begin{aligned}\hat{E} &= \hat{\sigma}_e^2 \mathbf{c}^\top (\mathbb{X}^\top \hat{\Sigma}^{-1} \mathbb{X})^{-1} \mathbf{c} \stackrel{\text{df}}{=} h(\hat{\sigma}_e^2, \hat{\boldsymbol{\theta}}), \\ \hat{V} &= \hat{\mathbf{g}}^\top \hat{\mathbb{I}}^{-1} \hat{\mathbf{g}},\end{aligned}$$

$\hat{\mathbb{I}}^{-1}$ is the estimated covariance matrix of $(\hat{\sigma}_e^2, \hat{\boldsymbol{\theta}}^\top)^\top$, and

$$\hat{\mathbf{g}} = \frac{\partial}{\partial(\sigma_e^2, \boldsymbol{\theta})} h(\hat{\sigma}_e^2, \hat{\boldsymbol{\theta}}).$$

6.5 Extended Linear Mixed Effects Model

6.5.1 Introduction

So far we have considered the model

$$\mathbf{Y}_i = \mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i + \boldsymbol{\varepsilon}_i, \quad i = 1, \dots, K,$$

where $\mathbf{b}_i \sim \mathbf{N}_q(\mathbf{0}, \mathbb{D})$ and $\boldsymbol{\varepsilon}_i \sim \mathbf{N}_{n_i}(\mathbf{0}, \sigma_e^2 \mathbb{I}_{n_i})$.

Now we remove the assumption of independence and homoskedasticity on the residual terms and assume instead

$$\boldsymbol{\varepsilon}_i \sim \mathbf{N}_{n_i}(\mathbf{0}, \sigma_e^2 \Lambda_i),$$

where Λ_i is $n_i \times n_i$ positive definite matrix parametrized by a vector of parameters $\boldsymbol{\lambda}$ of a fixed dimension. Marginally,

$$\mathbf{Y}_i \sim \mathbf{N}_{n_i}(\mathbb{X}_i \boldsymbol{\beta}, \mathbb{Z}_i \mathbb{D} \mathbb{Z}_i^\top + \sigma_e^2 \Lambda_i).$$

Thus, a part of the variance and covariance structure is ascribed to the random effects \mathbf{b}_i , another part is ascribed to the residual terms $\boldsymbol{\varepsilon}_i$. For a given covariance structure, there may be multiple ways to partition it between the random effects and residual terms, so one has to be careful with the specification of this model to make sure that all the parameters are identifiable.

6.5.2 Parameter estimation

There exists an invertible square root of Λ_i such that

$$\Lambda_i = (\Lambda_i^{1/2})^\top \Lambda_i^{1/2} \quad \text{and} \quad \Lambda_i^{-1} = \Lambda_i^{-1/2} (\Lambda_i^{-1/2})^\top.$$

Transform the observations, regressors and residual terms

$$\mathbf{Y}_i^* = (\Lambda_i^{-1/2})^\top \mathbf{Y}_i, \quad \mathbb{X}_i^* = (\Lambda_i^{-1/2})^\top \mathbb{X}_i, \quad \mathbb{Z}_i^* = (\Lambda_i^{-1/2})^\top \mathbb{Z}_i, \quad \boldsymbol{\varepsilon}_i^* = (\Lambda_i^{-1/2})^\top \boldsymbol{\varepsilon}_i.$$

Then $\text{var } \boldsymbol{\varepsilon}_i^* = \sigma_e^2 \mathbb{I}_{n_i}$ and \mathbf{Y}_i^* satisfy the standard one-level LME model

$$\mathbf{Y}_i^* = \mathbb{X}_i^* \boldsymbol{\beta} + \mathbb{Z}_i^* \mathbf{b}_i + \boldsymbol{\varepsilon}_i^*.$$

The Jacobian of the transformation is $|\Lambda_i|^{-1/2}$. The likelihood of the extended LME model is

$$L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2, \boldsymbol{\lambda} \mid \mathbf{Y}) = L(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma_e^2, \boldsymbol{\lambda} \mid \mathbf{Y}^*) \prod_{i=1}^K |\Lambda_i|^{-1/2}.$$

The first part is the likelihood of the standard LME model. The same holds for REML:

$$L_R(\boldsymbol{\theta}, \sigma_e^2, \boldsymbol{\lambda} \mid \mathbf{Y}) = L_R(\boldsymbol{\theta}, \sigma_e^2, \boldsymbol{\lambda} \mid \mathbf{Y}^*) \prod_{i=1}^K |\Lambda_i|^{-1/2}.$$

The parameters are estimated by the same decompositions of the log-likelihood as before, with an additional term $-1/2 \sum_{i=1}^K \log |\Lambda_i|$.

6.5.3 Generalized least squares

Consider a special case of the extended LME model with no random effects

$$\mathbf{Y} = \mathbb{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

where $\boldsymbol{\varepsilon} \sim \mathbf{N}_n(\mathbf{0}, \sigma_e^2 \Lambda)$. Suppose $\boldsymbol{\lambda}$ is known. Then

$$\mathbf{Y}^* = \mathbb{X}^* \boldsymbol{\beta} + \boldsymbol{\varepsilon}^*$$

satisfies the classical linear model. The LSE of $\boldsymbol{\beta}$ is $\widehat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) = [(\mathbb{X}^*)^\top \mathbb{X}^*]^{-1} (\mathbb{X}^*)^\top \mathbf{Y}^*$. The parameter σ_e^2 is estimated by the residual sum of squares

$$\widehat{\sigma}_e^2(\boldsymbol{\lambda}) = \frac{1}{n-p} \left\| \mathbf{Y}^* - \mathbb{X}^* \widehat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) \right\|^2.$$

(this is REML estimator). Plug this into the likelihood to get

$$\ell_M(\boldsymbol{\lambda}) = C - n \log \left\| \mathbf{Y}^* - \mathbb{X}^* \widehat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) \right\| - \frac{1}{2} \log |\Lambda|$$

and maximize it over $\boldsymbol{\lambda}$ to get the MLE. Or use the REML log-likelihood

$$\ell_R(\boldsymbol{\lambda}) = C - (n-p) \log \left\| \mathbf{Y}^* - \mathbb{X}^* \widehat{\boldsymbol{\beta}}(\boldsymbol{\lambda}) \right\| - \frac{1}{2} \log |\Lambda| - \frac{1}{2} \log \left| (\mathbb{X}^*)^\top \mathbb{X}^* \right|$$

and maximize it over $\boldsymbol{\lambda}$ to get the REML estimator.

This method is called *generalized least squares*.

6.5.4 Decomposing variance structure

The ability to specify arbitrary Λ_i (subject to identifiability assumption) allows to include in the model correlated residual terms with time series or spatial structure and/or to model unequal variances of residual terms.

Decompose

$$\varepsilon_{ij} = W_{ij} + \eta_{ij},$$

where η_{ij} are independent variables with equal variance σ_e^2 and W_{ij} are values of a random Gaussian processes with zero mean, independent of η_{ij} . We can take, for example:

- $W_{ij} = W_i(j)$ random processes in discrete time ($j = 1, \dots, n_i$ are indices of ordered observations within subject), independent between subjects.
- $W_{ij} = W_i(t_{ij})$ random processes in continuous time (t_{ij} is the time of the observation Y_{ij}), independent between subjects.
- $W_{ij} = W_i(u_{ij}, v_{ij})$ random fields on the plane ($[u_{ij}, v_{ij}]$ are the coordinates of the observation Y_{ij}), independent between subjects.

Consider $W_{ij} = W_i(t_{ij})$ as an example. Let $W_i(t)$ be weakly stationary with $\mathbf{E} W_i(t) = 0$, $\text{var } W_i(t) = \tau^2$, $\text{cor}(W_i(t), W_i(s)) = \rho(|t - s|)$. Then the variance is decomposed into three components

$$\text{var } Y_{ij} = \mathbf{Z}_{ij}^T \mathbb{D} \mathbf{Z}_{ij} + \tau^2 + \sigma_e^2,$$

the random effects variance $\mathbf{Z}_{ij}^T \mathbb{D} \mathbf{Z}_{ij}$ that depends on regressors \mathbf{Z}_{ij} , the variance of the serial component τ^2 and the white noise variance σ_e^2 . The covariance has two components

$$\text{cov}(Y_{ij}, Y_{ik}) = \mathbf{Z}_{ij}^T \mathbb{D} \mathbf{Z}_{ik} + \tau^2 \rho(|t_{ij} - t_{ik}|).$$

Denote by \mathbb{R}_i the correlation matrix of $(W_i(t_{i1}), \dots, W_i(t_{in_i}))$. We get

$$\text{var } \mathbf{Y}_i = \mathbf{Z}_i \mathbb{D} \mathbf{Z}_i^T + \tau^2 \mathbb{R}_i + \sigma_e^2 \mathbb{I}_{n_i}$$

and

$$\Lambda_i = \frac{\tau^2}{\sigma_e^2} \mathbb{R}_i + \mathbb{I}_{n_i}.$$

The process $W_i(t)$ is chosen to generate the desired autocorrelation function, for example

- exponential correlation $\rho(u) = \exp\{-\lambda u\}$
- Gaussian correlation $\rho(u) = \exp\{-\lambda u^2\}$
- compound symmetry $\rho(u) = \lambda$ (beware nonidentifiability of the random intercept)
- $AR(1)$ process $W_{ij} = \lambda W_{i,j-1} + \nu_{ij}$, with correlations $\rho(u) = \lambda^{|j-k|}$
- $ARMA(p, q)$ process

For spatial correlations, one can take some distance function d (Euclidean, L_1 , maximal) and calculate distances between measurements $d_{ijk} = d([u_{ij}, v_{ij}], [u_{ik}, v_{ik}])$. Then specify $\rho(d)$, the correlation between two measurements taken at the distance d of each other, for example

- exponential correlation $\rho(d) = \exp\{-\lambda d\}$
- Gaussian correlation $\rho(d) = \exp\{-\lambda d^2\}$
- linear correlation $\rho(d) = (1 - \lambda d)I(d < 1/\lambda)$

There are graphical methods for the choice of the right correlation structure (sample autocorrelation function, sample variogram). For more details, see [Diggle et al. \(2002, Chapter 5\)](#).

6.6 Comparison of LME and GEE Approaches

Group-dependent data with linear mean structure can be analyzed either by LME or by GEE methods. Both methods provide consistent estimators of mean effects β as long as the mean structure is correct and the number of independent groups K grows to infinity.

However, LME and GEE differ in some important aspects:

LME

- specifies a detailed model for $\text{var } \mathbf{Y}_i$
- assumes normality of \mathbf{b}_i and ε_i
- yields important information on the variance structure (decomposition of variance, estimates of variance components, random effects, hypotheses tests about variance structure)
- inference on β (tests, confidence intervals) is invalid if the variance structure is not correctly specified
- needs relatively large K , its performance when K is moderate or small is uncertain

GEE

- uses a working model for $\text{var } \mathbf{Y}_i$, which is not assumed to be correct
- does not make assumptions on the distribution of \mathbf{Y}_i
- does not provide much information on the variance structure
- inference on β is valid even if the working variance structure is incorrect, as long as K is large enough
- needs relatively large K , fails when K is small

7 Generalized Linear Mixed Models

7.1 Model and Assumptions

Again, we observe K independent random vectors $\mathbf{Y}_1, \dots, \mathbf{Y}_K$, where

$$\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^\top, \quad i = 1, \dots, K.$$

In this chapter, we develop a model for non-normal group-dependent data by introducing random effects into the generalized linear model.

Definition 7.1. The data \mathbf{Y}_i satisfy the *generalized linear mixed model** [GLMM] if

1. $\mathbf{Y}_1, \dots, \mathbf{Y}_K$ are independent.
2. There exist unobserved iid q -dimensional random vectors \mathbf{b}_i with density $h(\mathbf{b}; \boldsymbol{\psi})$ parametrized by an s -dimensional parameter $\boldsymbol{\psi}$ such that Y_{i1}, \dots, Y_{in_i} are conditionally independent given \mathbf{b}_i and the conditional density of Y_{ij} given \mathbf{b}_i has the form of the GLM

$$f(y | \mathbf{b}_i) = \exp\left\{\frac{y\theta_{ij} - b(\theta_{ij})}{\varphi} + c(y, \varphi)\right\}.$$

3. The canonical parameter θ_{ij} depends on the p -dimensional fixed effects covariates \mathbf{X}_{ij} , fixed regression parameters $\boldsymbol{\beta} \in \mathbb{R}^p$, random effects covariates \mathbf{Z}_{ij} of dimension $q \leq p$, and random effects \mathbf{b}_i through the linear predictor

$$\eta_{ij} = \mathbf{X}_{ij}^\top \boldsymbol{\beta} + \mathbf{Z}_{ij}^\top \mathbf{b}_i.$$

4. There exists a link function g such that $g(\mu_{ij}) = \eta_{ij}$, where $\mu_{ij} = b'(\theta_{ij}) \equiv \mathbb{E}[Y_{ij} | \mathbf{b}_i]$.

Note.

- Conditionally on \mathbf{b}_i , the response Y_{ij} satisfies the GLM. The presence of the common \mathbf{b}_i in all Y_{i1}, \dots, Y_{in_i} brings in within-group correlations.
- The random effects covariates \mathbf{Z}_{ij} are usually taken as a subset of the fixed effects covariates \mathbf{X}_{ij} .

* Český zobecněný lineární smíšený model

- Commonly, the random effects are assumed to be normally distributed

$$\mathbf{b}_i \sim \mathbf{N}_q(\mathbf{0}, \mathbb{D}),$$

where the covariance matrix \mathbb{D} is parametrized by a vector of parameters $\boldsymbol{\psi}$. Then $h(\mathbf{b}; \boldsymbol{\psi})$ is a multivariate normal density.

Note. The conditional moments of Y_{ij} given \mathbf{b}_i are

$$\begin{aligned} \mathbb{E}[Y_{ij} | \mathbf{b}_i] &= g^{-1}(\mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i) \\ \text{var}[Y_{ij} | \mathbf{b}_i] &= \varphi V(\mu_{ij}), \end{aligned}$$

where $V(\mu) = b''(\theta)$ is the variance function. The expectation of Y_{ij} given the covariates is

$$\mathbb{E} Y_{ij} = \mathbb{E} \mathbb{E}[Y_{ij} | \mathbf{b}_i] = \mathbb{E} g^{-1}(\mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i),$$

where the expectation on the right-hand side is over the distribution of the random effects (conditionally on all covariates). Unless g is linear, the distribution of Y_{ij} given the covariates does not follow the generalized linear model.

7.2 Parameter Estimation

The likelihood for data that satisfy Definition 7.1 can be written as

$$L(\boldsymbol{\beta}, \boldsymbol{\psi} | \mathbf{Y}) = \prod_{i=1}^K \int \prod_{j=1}^{n_i} f(Y_{ij} | \boldsymbol{\beta}, \mathbf{b}_i) h(\mathbf{b}_i | \boldsymbol{\psi}) d\mathbf{b}_i.$$

Unlike in the linear mixed effects model, here the integral cannot be calculated explicitly even if h is normal.

If the link g is canonical and h is multivariate normal, the likelihood can be expressed as

$$\begin{aligned} L(\boldsymbol{\beta}, \boldsymbol{\psi} | \mathbf{Y}) &= \prod_{i=1}^K (2\pi)^{-q/2} |\mathbb{D}|^{-1/2} \times \\ &\times \int \exp \left\{ \frac{1}{\varphi} \left[\mathbf{Y}_i^T (\mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i) - \mathbf{1}_{n_i}^T b(\mathbb{X}_i \boldsymbol{\beta} + \mathbb{Z}_i \mathbf{b}_i) \right] + \mathbf{1}_{n_i}^T c(\mathbf{Y}_i, \varphi) - \frac{1}{2} \mathbf{b}_i^T \mathbb{D}^{-1} \mathbf{b}_i \right\} d\mathbf{b}_i, \end{aligned}$$

where the functions b and c are applied to vectors by element-by-element calculation. The log-likelihood has the form

$$\ell(\boldsymbol{\beta}, \boldsymbol{\psi} | \mathbf{Y}) = -\frac{K}{2} \log |\mathbb{D}| + \sum_{i=1}^K \log \int \exp \{ \dots \} d\mathbf{b}_i + C.$$

Many different approaches have been suggested to maximize $\ell(\boldsymbol{\beta}, \boldsymbol{\psi} \mid \mathbf{Y})$, the use of the Laplace approximation is one among them.

The Laplace approximation is

$$\int_{\mathbb{R}^q} e^{Q(\mathbf{b})} d\mathbf{b} \approx (2\pi)^{q/2} \left| -Q''(\tilde{\mathbf{b}}) \right|^{-1/2} \exp\{Q(\tilde{\mathbf{b}})\},$$

where $\tilde{\mathbf{b}}$ is the point where the maximum of the function Q is attained. This approximation is obtained by replacing Q in the integrand by the second order Taylor expansion of Q around $\tilde{\mathbf{b}}$ and integrating the exponentiated quadratic function as a Gaussian density.

After plugging the Laplace approximation into the log-likelihood, the integral disappears and the resulting expression can be maximized using a combination of a modified IWLS algorithm for estimating $\boldsymbol{\beta}$, modified Henderson's equations for estimating \mathbf{b}_i , and moment estimation for $\boldsymbol{\psi}$ and φ . There are many different approaches to implementing these ideas. The calculations that need to be performed and the development of formulas and algorithms are rather tedious.

One of approaches leads to the following estimation procedure. At each iterative step, using the parameter estimates obtained at the previous step, calculate successively:

1. The IWLS step:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \widehat{\mathbb{W}}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \widehat{\mathbb{W}}^{-1} \widehat{\mathbf{Y}}^*,$$

where

$$\widehat{\mathbb{W}} = \begin{pmatrix} \widehat{\mathbb{W}}_1 & \cdots & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \cdots & \widehat{\mathbb{W}}_K \end{pmatrix}, \quad \widehat{\mathbb{W}}_i = \text{diag} \{ \widehat{\varphi} V(\widehat{\mu}_{ij}) [g'(\widehat{\mu}_{ij})]^2 \}_{j=1}^{n_i} + \mathbf{Z}_i \widehat{\mathbb{D}} \mathbf{Z}_i^T$$

$$\text{and } \widehat{\mathbf{Y}}_{ij}^* = g(\widehat{\mu}_{ij}) + (Y_{ij} - \widehat{\mu}_{ij}) g'(\widehat{\mu}_{ij}).$$

2. The Henderson step:

$$\widehat{\mathbf{b}}_i = \widehat{\mathbb{D}} \mathbf{Z}_i^T \widehat{\mathbb{W}}_i^{-1} (\widehat{\mathbf{Y}}_i^* - \mathbf{X}_i \widehat{\boldsymbol{\beta}}).$$

3. Moment estimation of an unstructured \mathbb{D} :

$$\begin{aligned} \widehat{\mathbb{D}} &= \frac{1}{K} \sum_{i=1}^K \widehat{\mathbf{E}}(\mathbf{b}_i \mathbf{b}_i^T) = \frac{1}{K} \sum_{i=1}^K [\widehat{\mathbf{E}}(\mathbf{b}_i \mid \mathbf{Y}_i)]^{\otimes 2} + \frac{1}{K} \sum_{i=1}^K \widehat{\text{var}}(\mathbf{b}_i \mid \mathbf{Y}_i) = \\ &= \frac{1}{K} \sum_{i=1}^K \widehat{\mathbf{b}}_i^{\otimes 2} + \frac{1}{K} \sum_{i=1}^K \widehat{\text{var}}(\mathbf{b}_i \mid \widehat{\mathbf{b}}_i). \end{aligned}$$

The last term is calculated according to Proposition 6.2.

7.3 Interpretation

Recall that

$$\mathbb{E}[Y_{ij} | \mathbf{b}_i] = g^{-1}(\mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i) \quad (7.1)$$

and

$$\mathbb{E} Y_{ij} = \mathbb{E} g^{-1}(\mathbf{X}_{ij}^T \boldsymbol{\beta} + \mathbf{Z}_{ij}^T \mathbf{b}_i). \quad (7.2)$$

Thus, $\boldsymbol{\beta}$ in (7.1) expresses the effect of \mathbf{X}_{ij} on $\mathbb{E} Y_{ij}$ conditionally on \mathbf{b}_i , that is, when the latent characteristics of the subject do not change. These effects are called *the subject-specific effects*. They describe the effects on $\mathbb{E}[Y_{ij} | \mathbf{b}_i]$ when the given subject changes the value of \mathbf{X}_{ij} . In general, the parameters $\boldsymbol{\beta}$ do not compare two *different* subjects that differ in the value of \mathbf{X}_{ij} .

The effect of \mathbf{X}_{ij} on $\mathbb{E} Y_{ij}$ unconditionally on \mathbf{b}_i , that is, in the general population, is called *the population-averaged effect*. Because $\mathbb{E} Y_{ij}$ in (7.2) does not necessarily have the form of a GLM with the link function g , the parameters $\boldsymbol{\beta}$ from (7.1) in general do not possess population-averaged interpretation.

There are several special cases when a relationship between the subject-specific and population-averaged model can be found:

- If the link function g is linear (including identity link as a special case) then the subject-specific and population-averaged model have the same form and the parameters $\boldsymbol{\beta}$ are the same in both. This is the case of the LME model of the previous chapter.
- If the link function g is logarithm and \mathbf{b}_i are normal then the subject-specific and population-averaged model have the same loglinear form and the parameters $\boldsymbol{\beta}$ are the same in both as long as the covariate does not appear in the random effects structure.
- If the link function g is probit and \mathbf{b}_i are normal then the subject-specific and population-averaged model have the same probit form. The parameters $\boldsymbol{\beta}$ in the subject-specific and population averaged model differ.

When the parameters in the subject-specific and population averaged model differ, their value in the population averaged model is always closer to zero than in the subject-specific model.

As an example, take binary Y_{ij} and logistic regression with a single continuous covariate X_{ij} and a random intercept. Let $\pi(x_{ij})$ be the conditional probability of success given x_{ij} and b_i . The model equation is

$$\log \frac{\pi(x_{ij})}{1 - \pi(x_{ij})} = \beta_0 + \beta_1 x_{ij} + b_i,$$

which can be rewritten as

$$\pi(x_{ij}) = \frac{\exp\{\beta_0 + \beta_1 x_{ij} + b_i\}}{1 + \exp\{\beta_0 + \beta_1 x_{ij} + b_i\}}.$$

The random intercept b_i can be interpreted as an innate propensity for success of the i -th subject that affects the odds of success at all trials performed by this subject. The parameter β_1 is the subject-specific effect of the covariate X_{ij} ; it explains what happens to the success probability of a given subject when the same subject is tested in different conditions.

The population-averaged model is

$$\pi^*(x_{ij}) = \mathbf{E} \frac{\exp\{\beta_0 + \beta_1 x_{ij} + b_i\}}{1 + \exp\{\beta_0 + \beta_1 x_{ij} + b_i\}},$$

where the expectation is taken over the distribution of b_i in the whole population. This is no longer a logistic regression model. However, if b_i is normally distributed with a small variance, $\pi^*(x_{ij})$ can be approximated by

$$\pi^*(x_{ij}) \approx \frac{\exp\{\beta_0^* + \beta_1^* x_{ij}\}}{1 + \exp\{\beta_0^* + \beta_1^* x_{ij}\}},$$

where $|\beta_1^*| < |\beta_1|$. The parameter β_1^* is the population-averaged effect of the covariate X_{ij} ; it explains what happens to the mean success probability of the whole population when the conditions change. The population-averaged effects are always weaker than the subject-specific effects.

7.4 Comparison of GLMM vs. GEE models

GLMM is a subject-specific model, its parameters have subject-specific interpretation. GEE is a population-averaged model, its parameters have population-averaged interpretation. In general, GLMM estimates different parameters than GEE. Also, GLMM and GEE models with the same link function cannot be both correct (except in special cases, some of which were mentioned in the previous section).

Otherwise, the main points of the discussion in Section 6.6 are still valid.

The choice between the GLMM and the GEE model should be driven mainly by the desired interpretation of the parameters.

Bibliography

- Diggle, P., Heagerty, P., Liang, K. and Zeger, S. (2002). *Analysis of Longitudinal Data*, Oxford Statistical Science Series, 2nd edn, Oxford University Press, Oxford.
- Harville, D. A. (1974). Bayesian inference for variance components using only error contrasts, *Biometrika* **61**(2): 383–385.
- Henderson, C. (1984). *Applications of linear models in animal breeding*, University of Guelph.
- Jennrich, R. I. and Schluchter, M. D. (1986). Unbalanced repeated-measures models with structured covariance matrices, *Biometrics* **42**(4): 805–820.
- Kenward, M. G. and Roger, J. H. (1997). Small sample inference for fixed effects from restricted maximum likelihood, *Biometrics* **53**(3): 983–997.
- Laird, N. M. and Ware, J. H. (1982). Random-effects models for longitudinal data, *Biometrics* **38**(4): 963.
- Liang, K. and Zeger, S. (1986). Longitudinal data analysis using generalized linear models, *Biometrika* **73**(1): 13–22.
- Lindstrom, M. J. and Bates, D. M. (1988). Newton-Raphson and EM algorithms for linear mixed-effects models for repeated-measures data, *Journal of the American Statistical Association* **83**(404): 1014–1022.
- Nelder, J. and Wedderburn, R. (1972). Generalized linear models, *Journal of the Royal Statistical Society. Series A* **135**(3): 370–384.
- Palmgren, J. (1981). The Fisher information matrix for log linear-models arguing conditionally on observed explanatory variables, *Biometrika* **68**(2): 563–566.
- Patterson, H. D. and Thompson, R. (1971). Recovery of inter-block information when block sizes are unequal, *Biometrika* **58**(3): 545–554.
- Prentice, R. L. (1988). Correlated binary regression with covariates specific to each binary observation, *Biometrics* **44**(4): 1033–1048.
- Prentice, R. L. and Zhao, L. P. (1991). Estimating equations for parameters in means and covariances of multivariate discrete and continuous responses, *Biometrics* **47**(3): 825–839.
- Satterthwaite, F. E. (1946). An approximate distribution of estimates of variance components, *Biometrics Bulletin* **2**(6): 110–114.

Bibliography

- Stram, D. O. and Lee, J. W. (1994). Variance components testing in the longitudinal mixed effects model, *Biometrics* **50**(4): 1171–1177.
- Stram, D. O. and Lee, J. W. (1995). Corrections: Variance component testing in the longitudinal mixed effects model, *Biometrics* **51**(3): 1196–1196.
- Wedderburn, R. (1974). Quasi-likelihood functions, generalized linear-models, and Gauss-Newton method, *Biometrika* **61**(3): 439–447.
- White, H. (1980). A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity, *Econometrica* **48**(4): 817–838.
- White, H. (1982). Maximum-likelihood estimation of mis-specified models, *Econometrica* **50**(1): 1–25.
- Yan, J. and Fine, J. (2004). Estimating equations for association structures, *Statistics in Medicine* **23**(6): 859–874.

Index

- adjusted dependent variable, **22**, 31, 39, 40, 43
- autocorrelation function, 99
- canonical link function, **16**, 21–23, 42, 101
- canonical parameter, **11**, 16, 18, 100
 - estimation, 14
- cauchit link, **35**
- cauchit regression model, **35**
- complementary log-log link, **36**
- conditional independence, **52**, 54, 55
- conditional odds ratio, **52**, 53–56
- contingency tables, 45–59
- Cook’s distance, 31
- correlation structure, 99
 - working, **68**, 69–72
- deviance, **24**, 28, 30, 40, 44, 50
- deviance residuals, **30**
 - standardized, **30**
- deviance test, **28**, 50
- dispersion parameter, **11**, 16, 60, 61
 - estimation, 15, 23, 66, 69
- distribution
 - alternative, 12, 14, 18, 34
 - beta, 60
 - beta-binomial, 60
 - binomial, 34, 60
 - gamma, 12, 13, 17, 61
 - geometric, 61
 - inverse Gaussian, 12, 13, 17
 - multinomial, 46–48
 - negative binomial, 61
 - normal, 11, 13, 17
 - Poisson, 12, 14, 17, 42, 46–48, 60
 - poisson-gamma, 61
- exponential family, **11**, 16, 17
- F test, 8, 75, 76, 78, 95
- fitted values, **21**
- generalized estimating equations, **68**, 99, 104
- generalized least squares, 97
- generalized linear mixed model, **100**
- generalized linear model, **16**, 100
- Henderson’s equations, 84, 86, 102
- iterative weighted least squares, **22**, 39, 40, 43, 65, 69, 102
- Kenward-Roger approximation, 94
- Laplace approximation, 102
- leverages, 30
- likelihood ratio test, 26, 28, 29, 63, 66, 92, 93
- linear mixed effects model, **81**
 - extended, 96–99
 - marginal form, 81, 83
 - marginal likelihood, 83
 - maximum likelihood, 86
 - multi-level, 82
 - REML, 87–92, 94, 97
 - single-level, 81
- linear predictor, **16**, 18, 31, 100
- link function
 - probit, **35**
- link function, **16**, 31, 61, 67, 100

- canonical, [16](#), [21–23](#), [42](#), [101](#)
- cauchit, [35](#)
- complementary log-log, [36](#)
- logistic, [18](#), [35](#)
- probit, [38](#)
- logistic link, [18](#), [35](#)
- logistic regression model, [18](#), [35](#), [38–42](#), [57–59](#)
- loglinear model, [17](#), [42–59](#)
- marginal independence, [51](#), [54](#)
- marginal likelihood in LME model, [83](#)
- marginal odds ratio, [51](#), [53–55](#)
- model
 - cauchit, [35](#)
 - logistic, [18](#), [35](#), [38–42](#), [57–59](#)
 - loglinear, [17](#), [42–59](#)
 - null, [19](#)
 - probit, [35](#), [38](#)
 - saturated, [19](#), [24](#), [49](#), [56](#)
- null model, [19](#)
- odds, [38](#), [49](#)
- odds ratio, [39](#), [50](#)
 - conditional, [52](#), [53–56](#)
 - marginal, [51](#), [53–55](#)
- offset, [44](#)
- overdispersion, [60](#), [61](#)
- parameter
 - canonical, [11](#), [16](#), [18](#), [100](#)
 - estimation, [14](#)
 - dispersion, [11](#), [16](#), [60](#), [61](#)
 - estimation, [15](#), [23](#), [66](#), [69](#)
- Pearson chi-square statistic, [24](#), [29](#), [41](#), [42](#), [44](#), [50](#), [66](#), [69](#)
- Pearson residuals, [29](#), [70](#)
 - standardized, [30](#)
- Poisson process, [44](#)
- population-averaged effect, [103](#)
- probit link, [35](#), [38](#)
- probit regression model, [35](#), [38](#)
- pseudo-score, [65](#), [68](#)
- quasi-likelihood, [62](#)
- quasi-score, [62](#)
- random effects, [75](#), [77](#), [81](#), [100](#)
- relative precision factor, [81](#), [86](#)
- REML, [87–92](#), [94](#), [97](#)
- residual maximum likelihood, *see* REML
- residuals
 - deviance, [30](#)
 - standardized, [30](#)
 - Pearson, [29](#), [70](#)
 - standardized, [30](#)
- restricted maximum likelihood, *see* REML
- sandwich variance estimator, [9](#), [64](#), [65](#), [66](#), [69](#)
- Satterthwaite’s approximation, [79](#), [94](#), [95](#)
- saturated model, [19](#), [24](#), [49](#), [56](#)
- subject-specific effect, [103](#)
- t test, [8](#), [94](#)
- variance components, [75](#), [77](#)
- variance function, [13](#), [61](#), [101](#)
 - working, [65](#), [66](#), [68](#)
- White estimator, *see* sandwich variance estimator
- working correlation matrix, [68](#), [69–72](#)
- working variance function, [65](#), [66](#), [68](#)