

Course overview

LM (TMA 4267)

Linear models, general linear model

GLM

Generalized linear models

Binomial, binomial

Poisson

Gamma

Multinomial (VGLM)

LMM

Linear mixed models, linear models  
including random effects, dependent data

REML

GLLMM

Generalized linear mixed models

Laplace approximation of marginal likelihood

REML

Introduction, binary regression

Recall

Multiple regression:

Assumptions: For  $i=1, 2, \dots, n$ 

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_k X_{ik} + e_i$$

where  $e_i \stackrel{iid}{\sim} N(0, \sigma^2)$ .

In matrix notation

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{\varepsilon}, \quad \underline{\varepsilon} \sim N(0, \sigma^2 I_n).$$

MLE of  $\underline{\beta}$

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

Binary regression example:  $\underline{Y}_i \in \{0, 1\}$

Aim: How does a binary outcome depend on covariates of interest  $x_{i1}, x_{i2}, \dots, x_{ik}$ .

Model: For  $i=1, 2, \dots, n$   
indep.

$$Y_i \sim \text{Bernoulli}(p_i),$$

where

$$\ln \frac{p_i}{1-p_i} = \underbrace{\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}}_{\text{linear predictor } \eta_i} + \underbrace{\epsilon_i}_{\substack{\text{no error} \\ \text{term and} \\ \text{dispersion param.}}}$$

$= \text{logit } p_i$

Note: The logit function  $\text{logit } p = \ln \frac{p}{1-p}$  maps  $p \in (0, 1)$  to  $\mathbb{R}$ .

The inverse:

$$\ln \frac{p}{1-p} = \eta$$

$$p = (1-p)e^\eta$$

$$(1 + e^\gamma) p = e^\gamma$$

$$p = \frac{e^\gamma}{1 + e^\gamma} = \frac{1}{1 + e^{-\gamma}} = \text{logit}^{-1}(\gamma)$$

Thus

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik})}}$$

Odds interpretation :-

$$\frac{p_i}{1-p_i} = \text{Odds}(X_i = 1).$$

The model assumes that

$$\ln \frac{p_i}{1-p_i} = \beta_0 + \beta_1 x_{i1} + \dots$$

$\underbrace{\quad}_{\text{odds}}$

} additive effects on log odds of event

$$\frac{p_i}{1-p_i} = e^{\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}}$$

↑

} multiplicative effects on odds of event

One unit change in e.g.  $x_{i1}$  changes the odds by a factor (an odds ratio) of  $e^{\beta_1}$ .

Example:  $\beta_1 = 0.1 \Rightarrow$  odds ratio is

$$e^{\beta_1} = e^{0.1} \approx 1.11$$

i.e. the odds increase by 11%.

Parameter estimation:  
Maximise the likelihood (numerically)

$$L(\beta) = P(Y_1 = y_1 \cap \dots \cap Y_n = y_n)$$

$$= \prod_{i=1}^n P(Y_i = y_i)$$

$$= \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i}$$

where for  $i=1, \dots, n$

$$p_i = \text{logit}^{-1}(\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik})$$

### Exponential family

Def.: Pdf or pmf on the form

$$f(y|\theta) = \exp \left( \underbrace{\frac{y\theta - b(\theta)}{\phi}}_w + c(y, \phi, w) \right)$$

depends only  
on  $\theta$

don't  
depend on  
 $\theta$

and

$$\ln f(y|\theta) = \frac{y\theta - b(\theta)}{\phi} w + c(y, \phi, w)$$

where  $\theta$  is a parameter (the natural or canonical),  $b(\theta)$  is a twice

differentiable function,  $\phi$  is a dispersion parameter,  $\omega$  is a known constant.

Support not dependent on  $\Theta$ .

Ex.:  $Y \sim \text{Bernoulli}(\pi)$

$$\text{Pmf } f(y|\pi) = \pi^y (1-\pi)^{1-y}$$

$$\ln f(y|\pi) = y \ln \pi + (1-y) \ln (1-\pi)$$

$$= y \underbrace{[\ln \pi - \ln(1-\pi)]}_{\theta} + \ln(1-\pi)$$

$$= y \underbrace{\ln \frac{\pi}{1-\pi}}_{\theta} - \underbrace{(-\ln(1-\pi))}_{b(\theta)}$$

$$\theta = \log \pi$$

$$\pi = e^{\theta}$$

$$= \frac{1}{1+e^{-\theta}} \doteq \frac{e^\theta}{1+e^\theta}$$

$$1-\pi = \frac{e^{-\theta}}{1+e^{-\theta}} = \frac{1}{1+e^\theta}$$

$$= y\theta + \ln \frac{1}{1+e^\theta}$$

$$= y\theta - \ln(1+e^\theta)$$

that is,  $\theta = \ln \frac{\pi}{1-\pi}$  is the natural parameter

and  $b(\theta) = \ln(1+e^\theta)$ ,  $c(y, \omega, \theta) = \omega$

Theorem: If  $Y$  belongs to the exponential family.

$$EY = b'(\theta) \quad (*)$$

$$\text{Var} Y = \frac{\phi}{\psi} b''(\theta)$$

Ex. cont:  $b'(\theta) = \frac{e^\theta}{1 + e^\theta} = \pi = EY$

$$\frac{\phi}{\psi} b''(\theta) = \frac{e^\theta (1 + e^\theta) - e^\theta e^\theta}{(1 + e^\theta)^2}$$

$$= \frac{e^\theta}{(1 + e^\theta)^2} \cdot \frac{1}{1 + e^\theta} = \pi(1 - \pi)$$

$$= \text{Var} Y$$

Proof of (\*):  $Y$  has mgf

$$M_Y(u) = E e^{uY} = \int e^{\frac{y\theta - b(\theta)}{\phi} w + yu + c(y, \phi, w)} dy$$

$$= e^{\underbrace{c(y, \phi, w)}_{\text{const}} + (\theta + u)\frac{\phi}{\psi}}$$

$$+\frac{w}{\phi}b\left(\theta + \frac{\phi}{w}u\right) - \frac{w}{\phi}b(\theta)$$

$$= e^{\frac{w}{\phi}\left(b\left(\theta + \frac{\phi}{w}u\right) - b(\theta)\right)} = e^{\frac{w}{\phi}\left(\frac{\gamma(\theta + \frac{\phi}{w}u) - b(\theta + \frac{\phi}{w}u)}{\phi}\right)w + c(\gamma, \phi, w)}$$

Recall:

$$\begin{aligned} \frac{d}{du} M_X(u) &= \frac{d}{du} E(e^{uX}) \\ &= E \frac{d}{du} e^{uX} = E X e^{uX} \Big|_{u=0} = EX \end{aligned}$$

$$\frac{d^2}{du^2} M_X(u) \Big|_{u=0} = E(X^2)$$

$$\text{We want to find } \text{Var}X = E(X^2) - (Ex)^2$$

Def.: The cumulant generating function of a random variable  $X$  is

$$K_X(u) = \ln M_X(u).$$

$$\text{Lemma: } K_X'(u) = \frac{M'_X(u)}{M_X(u)} \Big|_{u=0} = \frac{M'_X(0)}{M_X(0)} = \frac{Ex}{1}$$

$$K_X^{''}(u) = \frac{M_X^{''}(u)M_X(u) - (M_X'(u))^2}{(M_X(u))^2}$$

$$\stackrel{u=0}{=} E(X^2) - (Ex)^2 = \text{Var}(X)$$

$$K^{''}(0) = \dots = E[(X-Ex)^3] \quad (\text{exercise 1b})$$

$$K^{(4)}(0) = \dots \neq E((X-Ex)^4) \quad (\text{see wikipedia})$$

$Y$  has cumulant generating

$$K_Y(u) = \frac{\omega}{\phi} \left( b\left(\theta + \frac{\phi}{\omega} u\right) - b(\theta) \right)$$

and

$$K_Y'(u) = \frac{\omega}{\phi} \left( b'\left(\theta + \frac{\phi}{\omega} u\right) \cdot \frac{\phi}{\omega} \right) \stackrel{u=0}{=} b'(\theta)$$

$$K_Y^{''}(u) = b''\left(\theta + \frac{\phi}{\omega} u\right) \frac{\phi}{\omega} \stackrel{u=0}{=} \frac{\phi}{\omega} b''(\theta)$$

which completes the proof.

## Exponential family cont.

31.8

Exponential distribution:

$$f(y|\lambda) = \lambda e^{-\lambda y}$$

$$= e^{-\lambda y + \ln \lambda}$$

$$= e^{y\theta - b(\theta)}$$

$$\text{so } \theta = -\lambda, b(\theta) = -\ln \lambda = \underline{-\ln(-\theta)} = -\frac{1}{\ln \theta}$$

$$EY = b'(\theta) = -\frac{1}{-\theta}(-1) = -\frac{1}{\theta} = \frac{1}{\lambda}$$

$$\text{Var } Y = \frac{\phi}{w} b''(\theta) = \frac{1}{\theta^2} = \frac{1}{(-\lambda)^2} = \frac{1}{\lambda^2}$$

Poisson:

$$f(y|\lambda) = \frac{e^{-\lambda} \lambda^y}{y!} = e^{y \ln \lambda - \lambda - \ln y!}$$

$$= e^{y\theta - b(\theta) - c(y, \phi, w)}$$

$$\text{where } \theta = \ln \lambda, b(\theta) = \lambda = e^\theta$$

$$EY = b'(\theta) = e^\theta = \lambda$$

$$\text{Var } Y = b''(\theta) = e^\theta = \lambda$$

Normal:

$$f(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

$$-\frac{\gamma^2}{2\sigma^2} + \frac{\gamma\mu}{\sigma^2} - \underbrace{\frac{\mu^2}{2\sigma^2}}_{c(\gamma, \phi, \omega)} - \frac{1}{2} \ln(2\pi\sigma^2)$$

$c(\gamma, \phi, \omega)$

$$= e^{\frac{\gamma\mu - \mu^2/2}{\sigma^2} - \frac{\gamma^2}{2\sigma^2} - \frac{1}{2} \ln(2\pi\sigma^2)}$$

$= e$

that is  $\theta = \mu$ ,  $b(\theta) = \underline{\mu^2/2} = \theta^2/2$ ,  $\phi = \sigma^2$

$$EY = b'(\theta) = \mu$$

$$\text{Var } Y = \frac{\phi}{\omega} b''(\theta) = \frac{\phi}{\omega} = \sigma^2$$

## Generalized linear models in general

For  $i=1, 2, \dots, n$  each observation  $y_i$  belongs to the exponential family and the mean

$$E(y_i) = \mu_i$$

is linked to a linear predictor

$$y_i = x_i^\top \beta, \quad x_i = \begin{pmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{ik} \end{pmatrix}$$

via a strictly monotonic, twice differentiable link function  $g$  or response function  $h = g^{-1}$  such that

$$\mu_i = h(y_i)$$

or

$$g(\mu_i) = y_i$$

If  $g(\mu_i)$  is the canonical parameter  
then  $g$  is the canonical link function.

## Multiple linear regression

$$Y = X\beta + \epsilon$$

$\downarrow \downarrow \downarrow$

$$n \times 1 \quad n \times p \quad p \times 1 \quad n \times 1$$

As a glm: For  $i = 1, 2, \dots, n$

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

where

$$\mu_i = \eta_i \quad (\text{identity link function})$$

and

$$\eta_i = \bar{x}_i^\top \beta$$

Parameter estimation: Maximise likelihood

$$L(\beta) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2} (y - X\beta)^\top (y - X\beta)}$$

or the log likelihood

$$l(\beta) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\gamma - X\beta)^T (\gamma - X\beta)$$

Score function: Want to solve the eq.

$$s(\beta) = \frac{\partial}{\partial \beta} l(\beta) = 0$$

### Matrix calculus

Def.: If  $y$  is a scalar and  $x$  is a  $p \times 1$  column vector, then

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \end{bmatrix}.$$

Rule 1: Thus, if  $a$  is another  $p \times 1$  column vector,

$$\frac{\partial}{\partial x} (a^T x) = \frac{\partial}{\partial x} (a_1 x_1 + \dots + a_p x_p)$$

$$= \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{bmatrix} = a.$$

Analogous to

$$\frac{d}{dx}(ax) = a. (x)$$

Def.: If  $\underline{Y}$  is a  $q \times 1$  column vector, then

$$\frac{\partial}{\partial \underline{x}} \underline{Y}^T = \begin{bmatrix} \frac{\partial Y_1}{\partial x_1} & \frac{\partial Y_2}{\partial x_1} & \dots \\ \frac{\partial Y_1}{\partial x_2} & & \\ \vdots & & \end{bmatrix}.$$

$p \times 1 \quad 1 \times q$

See wikipedia  
for other  
conventions

Rule 2: Thus, if  $A$  is  $q \times p$  matrix, then

$$\begin{aligned} \frac{\partial}{\partial \underline{x}} (A \underline{x})^T &= \frac{\partial}{\partial \underline{x}} \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1p}x_p & a_{21}x_1 + a_{22}x_2 + \dots \end{bmatrix} \\ &= \begin{bmatrix} a_{11} & a_{21} \\ a_{12} & \vdots \\ \vdots & \end{bmatrix} = A^T. \end{aligned}$$

Again, analogous to (\*)

If  $A$  is a  $p \times p$ -matrix, then

$$\begin{aligned} \frac{\partial}{\partial \underline{x}} (\underline{x}^T A \underline{x})^{p=2} &= \frac{\partial}{\partial \underline{x}} (a_{11}x_1^2 + a_{12}x_1x_2 + a_{21}x_1x_2 + a_{22}x_2^2) \\ &= \begin{bmatrix} 2a_{11}x_1 + a_{12}x_2 + a_{21}x_2 \\ a_{12}x_1 + a_{21}x_1 + 2a_{22}x_2 \end{bmatrix} \end{aligned}$$

$$= (\underline{A} + \underline{A}^T) \underline{x}.$$

Rule 3: If  $\underline{A}$  is symmetric

$$\frac{\partial}{\partial \underline{x}} (\underline{x}^T \underline{A} \underline{x}) = 2 \underline{A} \underline{x}.$$

Analogous to  $\frac{d}{dx} (ax^2) = 2ax$

Recall: If  $\underline{g}(\underline{x})$  is a  $q \times 1$  column-vector-valued function and  $f$  and  $x$  are scalars, then

$$\frac{\partial}{\partial \underline{x}} f(\underline{g}(\underline{x})) = \underbrace{\frac{\partial f}{\partial g_1} \cdot \frac{\partial g_1}{\partial \underline{x}} + \dots + \frac{\partial f}{\partial g_q} \cdot \frac{\partial g_q}{\partial \underline{x}}}_{\substack{1 \times q \\ 1 \times 1}} = \underbrace{\frac{\partial \underline{g}^T}{\partial \underline{x}} \cdot \frac{\partial f}{\partial \underline{g}}}_{\substack{q \times 1 \\ q \times 1}}$$

Rule 4 (chain rule):

$$\frac{\partial}{\partial \underline{x}} f(\underline{g}(\underline{x})) = \dots = \underbrace{\frac{\partial \underline{g}^T}{\partial \underline{x}}}_{\substack{P \times q \\ P \times P}} \cdot \underbrace{\frac{\partial f}{\partial \underline{g}}}_{\substack{q \times 1 \\ P \times 1}}$$

Multiple regression cont.:

$$s(\underline{\beta}) = \frac{\partial l}{\partial \underline{\beta}} = \frac{\partial}{\partial \underline{\beta}} \left( -\frac{1}{2\sigma^2} (\underline{y} - \underline{x}\underline{\beta})^T (\underline{y} - \underline{x}\underline{\beta}) \right)$$

$$= -\frac{1}{2\sigma^2} \cdot \frac{\partial}{\partial \underline{\beta}} f(\underline{g}(\underline{\beta})) \quad \underline{g}(\underline{\beta}) = \underline{y} - \underline{x}\underline{\beta}, \quad f(\underline{s}) = \underline{s}^T \underline{s}$$

$$= \frac{\partial g^T}{\partial \beta} \cdot \frac{\partial f}{\partial g}$$

$$= -\frac{1}{2\sigma^2} \cdot (-X^T) 2(Y - X\beta) \quad (\text{Rule 4, 3 and 2})$$

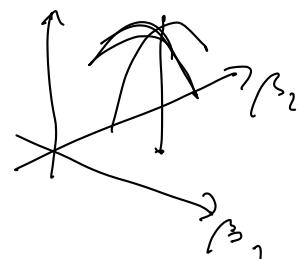
At the MLE of  $\beta$ ,

$$\hat{g}(\beta) = 0$$

$$X^T Y - X^T X \beta = 0$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$l(\beta)$$



Observed Fisher information

$$H(\beta) = - \underbrace{\frac{\partial^2}{\partial \beta \partial \beta^T}}_{p \times 1 \quad 1 \times p} l(\beta) = - \frac{\partial}{\partial \beta} \hat{g}^T(\beta)$$

$$= - \frac{\partial}{\partial \beta} \left( + \frac{1}{2\sigma^2} X^T (Y - X\beta) \right)^T = \frac{X^T X}{\sigma^2} \quad (\text{Rule 2})$$

Note how

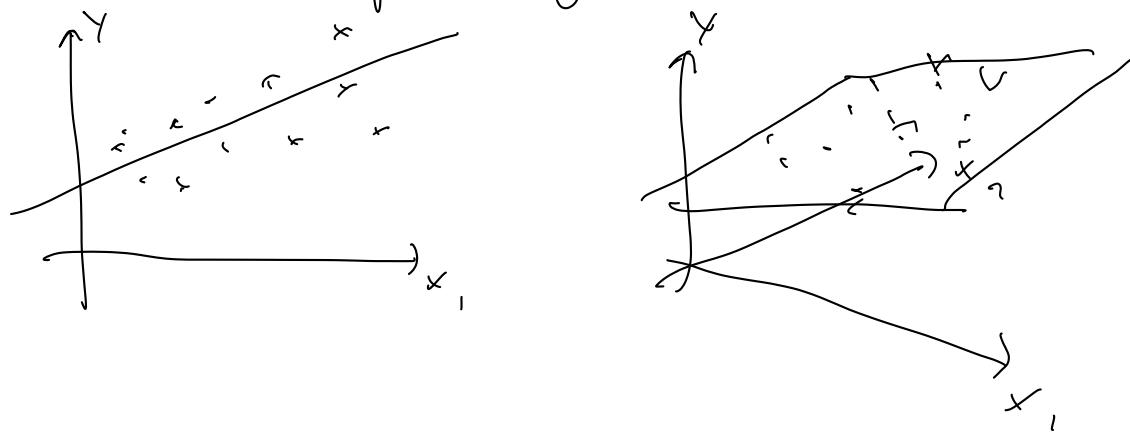
$$H^{-1}(\beta) = \sigma^2 (X^T X)^{-1} = \text{Var}(\hat{\beta})$$

In this case the (expected) Fisher information

$$F(\beta) \stackrel{\text{def}}{=} E\left(-\frac{\partial^2}{\partial \beta \partial \beta^T} l(\beta; Y)\right) = H(\beta),$$

the observed Fisher information.

Data centric, applied, geometric view:



Linear algebra, vector-space view:

$n$  dimensions

$$\underline{Y} \in \mathbb{R}^n$$

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{\epsilon}$$

$$= \begin{bmatrix} \underline{X}_1 & \underline{X}_2 & \dots & \underline{X}_p \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \underline{\epsilon}$$

$$= \underline{x}_1\beta_1 + \underline{x}_2\beta_2 + \dots + \underline{x}_p\beta_p + \underline{\epsilon}$$

We seek the linear combination

of the columns in  $\underline{X}$ ,  $\hat{\underline{y}} = \underline{X}\hat{\underline{\beta}} \in \text{colsp}(\underline{X})$   
 minimising  $\|\underline{y} - \hat{\underline{y}}\|^2 = (\underline{y} - \hat{\underline{y}})^T (\underline{y} - \hat{\underline{y}}) = (\underline{y} - \underline{X}\hat{\underline{\beta}})^T (\underline{y} - \underline{X}\hat{\underline{\beta}})$

The residual vector  $\hat{\underline{\epsilon}} = \underline{y} - \hat{\underline{y}}$  must therefore be orthogonal to all columns in  $\underline{X}$ . All the dot products  $\underline{x}_1^T \hat{\underline{\epsilon}}, \underline{x}_2^T \hat{\underline{\epsilon}}, \dots, \underline{x}_p^T \hat{\underline{\epsilon}}$  are therefore all zero, that is,

$$\left( \begin{array}{c} \underline{x}_1^T \\ \vdots \\ \underline{x}_p^T \end{array} \right)$$

$$\underline{X}^T \begin{pmatrix} \underline{\epsilon} \\ \hat{\beta} \end{pmatrix} = 0$$

$$\underline{X}^T (\underline{y} - \underline{X}\hat{\beta}) = 0$$

$$\hat{\beta} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} \quad (\text{the MLE})$$

The projection is given by

$$\hat{\underline{y}} = \underline{X}\hat{\beta} = \underbrace{\underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T}_{H} \underline{y} = H\underline{y}$$

$H$  has rank  $p$

The residuals is another projection

$$\hat{\underline{\epsilon}} = \underline{y} - \hat{\underline{y}} = \underline{y} - H\underline{y} = \underbrace{(\underline{I} - H)\underline{y}}$$

rank  $n-p$

Hence,

$$\left\| \frac{\hat{\underline{\epsilon}}}{\sigma} \right\|^2 = \frac{SSE}{\sigma^2} = \underbrace{\underbrace{z_1^2 + z_2^2 + \dots + z_{n-p}^2}_{n-p} + \underbrace{0^2 + 0^2 + \dots + 0^2}_p}_{n-p} \sim \chi^2_{n-p}$$

$z_1, z_2, \dots, z_n$  are  
new coordinates in a  
suitable orthonormal basis

(a rotation of the coordinate system)

and

$$E\left(\frac{SSE}{\sigma^2}\right) = n-p$$

such that

$\hat{\sigma}^2 = \frac{SSE}{n-p}$  is unbiased for  $\sigma^2$ .

$\hat{\epsilon}$  and  $\hat{\beta}$  are jointly multivariate normal with

$$\text{Cov}(\hat{\epsilon}, \hat{\beta}) = \text{Cov}\left(\underbrace{(\mathbb{I} - X(X^T X)^{-1} X^T) Y}_A, \underbrace{(X^T X)^{-1} X^T Y}_B\right)$$

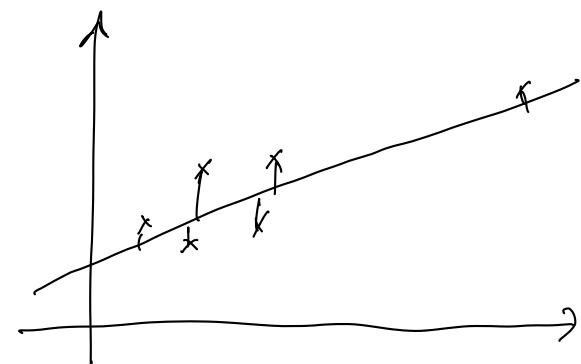
$$= A \text{Cov}(Y, Y) B^T$$

$$\begin{aligned} &= \sigma^2 \left( \mathbb{I} - X(X^T X)^{-1} X^T \right) X (X^T X)^{-1} \\ &= \sigma^2 \left( \underbrace{X(X^T X)^{-1}}_{\text{---}} - \underbrace{X(X^T X)^{-1} X^T X (X^T X)^{-1}}_{\cancel{X^T X (X^T X)^{-1}}} \right) \\ &= 0. \end{aligned}$$

Do the residuals  $\hat{\epsilon} = Y - \hat{Y}$  have constant variance? No:

$$\text{Var}(\hat{\epsilon}) = \text{Var}((\mathbb{I} - H) Y)$$

$$\begin{aligned} &= \sigma^2 (\mathbb{I} - H)(\mathbb{I} - H) \\ &= \sigma^2 (\mathbb{I} - 2H + \underline{H^2}) \\ &= \sigma^2 (\mathbb{I} - 2H + H) \\ &= \sigma^2 (\mathbb{I} - H) \end{aligned}$$



Standardized residuals

$$r_i = \frac{y_i - \hat{y}_i}{\sqrt{1 - h_{ii}}} \stackrel{\text{approx}}{\sim} N(0, 1)$$

If the model has an intercept, then since  $\hat{\epsilon}$  is orthogonal to all columns in  $X$  including the first column such that

$$\hat{\epsilon}^T \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = 0$$

$$\hat{\epsilon}_1 + \hat{\epsilon}_2 + \dots + \hat{\epsilon}_n = 0$$

T-tests:

$$\text{Var}(\hat{\beta}) = \sigma^2 (\hat{X}^T \hat{X})^{-1} = \sigma^2 \begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix}$$

It follows that

$$T = \frac{\hat{\beta}_j}{\sqrt{\frac{SSE}{\sigma^2} / (n-p)}} \sim t_{n-p}$$

$$\hat{\beta}_j \sim N(0, 1)$$

$$\sqrt{\frac{SSE}{\sigma^2} / (n-p)} \sim \chi_{n-p}^2$$

$$= \frac{\hat{\beta}_j}{\widehat{SE}(\hat{\beta}_j)}$$

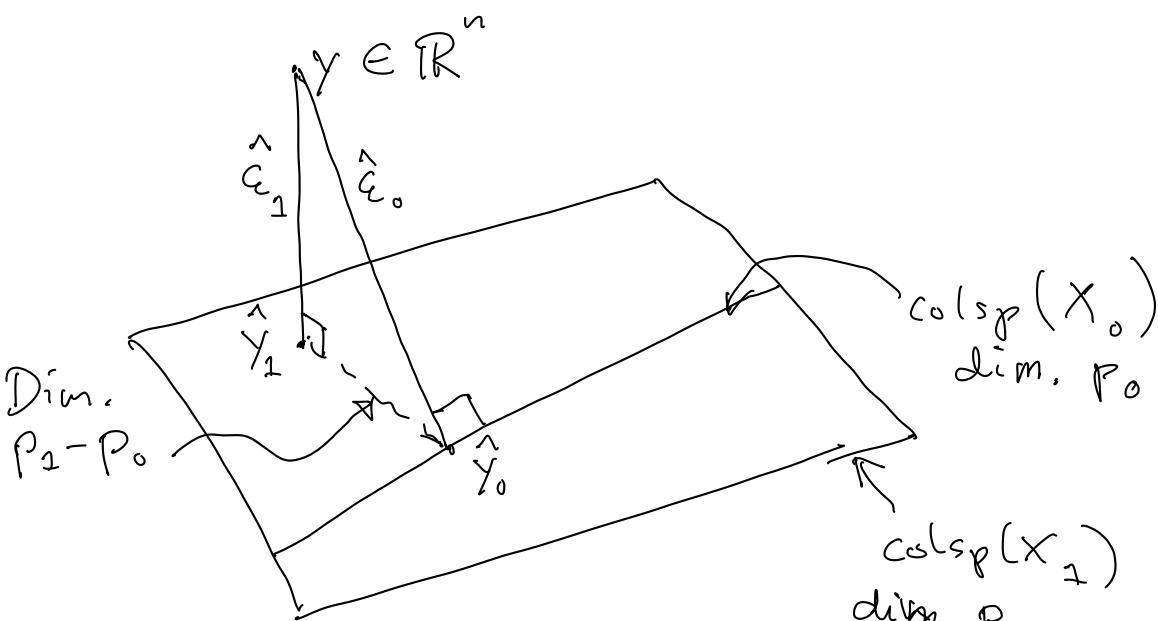
is t-distributed with  $n-p$  d.f.

under  $H_0: \beta_j = 0$

## The F-distribution

Testing:  $H_0: y = X_0 \beta_0 + \varepsilon$  vs  $H_1: y = X_1 \beta_1 + \varepsilon$

Vector-space view:



Orthogonality of  $\hat{\varepsilon}_0$  and  $\hat{y}_1 - \hat{y}_0$  implies that

$$\left\| \frac{\hat{y}_1 - \hat{y}_0}{\sigma} \right\|^2 \sim \chi^2_{P_1 - P_0} \text{ and indep. of } \sum \hat{\varepsilon}_1^2$$

under  $H_0$ . From the Pythagorean theorem we have

$$\left\| \hat{\varepsilon}_0 \right\|^2 = \left\| \hat{\varepsilon}_1 \right\|^2 + \left\| \hat{y}_1 - \hat{y}_0 \right\|^2$$

such that

$$\begin{aligned} \left\| \hat{y}_1 - \hat{y}_0 \right\|^2 &= \left\| \hat{\varepsilon}_0 \right\|^2 - \left\| \hat{\varepsilon}_1 \right\|^2 \\ &= SSE_0 - SSE_1. \end{aligned}$$

Hence,

$$\frac{SSE_0 - SSE_1}{\sigma^2} / (P_1 - P_0)$$

$$F = \frac{\text{SSE}_0 - \text{SSE}_1}{\sigma^2} / (P_1 - P_0)$$

$$\sim F_{P_1 - P_0, n - P_1}$$

$$\frac{SSE_1}{\sigma^2} / (n - p_1)$$

Reject  $H_0$  for large values of  $F$ .

Equivalent to a LRT (likelihood ratio test)  
(exercise 2b)

### Testing linear hypothesis

T-test: Testing  $H_0: \beta_i = \beta_{i0}$   
vs  $H_1: \beta_i \neq \beta_{i0}$ .

$$T = \frac{\hat{\beta}_i - \beta_{i0}}{SE(\hat{\beta}_i)} \sim t_{n-p}$$

Note that

$$T^2 = \frac{(\hat{\beta}_i - \beta_{i0})^2}{Var(\hat{\beta}_i)} \sim F_{1, n-p}$$

General linear hypothesis:

$$H_0: \underbrace{\sum_{r \times p} \beta_r}_{p \times 1} = \underbrace{d}_{r \times 1}$$

vs

$$H_1: C\beta \neq d$$

Under  $H_0$   $C\hat{\beta} - \underline{d}$  is normally distributed with

$$E(C\hat{\beta} - \underline{d}) = \underline{0}$$

and variance matrix

$$\begin{aligned} \text{Var}(C\hat{\beta} - \underline{d}) &= C \text{Var}(\hat{\beta}) C^T \\ &= \sigma^2 \underbrace{C(X^T X)^{-1} C^T}_{\Sigma} \end{aligned}$$

Hence

$$\underline{z} = \Sigma^{-1/2} (C\hat{\beta} - \underline{d})$$

has mean zero and

$$\text{Var}(\underline{z}) = \Sigma^{-1/2} \Sigma (\Sigma^{-1/2})^T = I$$

that is  $\underline{z} \sim N(0, I_n)$

Thus,

$$\|\underline{z}\|^2 = \underline{z}^T \underline{z} = z_1^2 + z_2^2 + \dots + z_n^2$$

$$= (C\hat{\beta} - \underline{d})^T \Sigma^{-1} (C\hat{\beta} - \underline{d})$$

$$= \frac{1}{\sigma^2} (C\hat{\beta} - \underline{d})^T (C(X^T X)^{-1} C^T)^{-1} (C\hat{\beta} - \underline{d})$$

$$\sim \chi_n^2$$

and

$$F = \frac{\frac{1}{n} (\hat{C}\beta - \underline{\alpha})^\top (C(X^\top X)^{-1} C) (\hat{C}\beta - \underline{\alpha}) / r}{\frac{SSE_1}{n-p} / (n-p)} \quad \text{Wald-test}$$
$$\sim F_{r, n-p}$$

How to fit  $H_0$ ? (Next lecture)

Fitting models with linear equality constraints

August 31

Fitting  $y = X\beta + \varepsilon$  under the constraint

$$C\beta = \underline{\alpha}$$

Fahrmeir p. 172-175: Lagrange:

$$\hat{\beta}^{(R)} = \underbrace{\hat{\beta}}_{\text{unrestricted MLE}} - \underbrace{(X^\top X)^{-1} C^\top (C(X^\top X)^{-1} C)^\top (\hat{C}\beta - \underline{\alpha})}_{\Delta_{H_0}}$$

It follows that

$$SSE_0 = \|\hat{\varepsilon}_0\|^2 = \|\hat{\varepsilon}_0 + X\Delta_{H_0}\|^2$$

$$= \|\hat{\varepsilon}_1\|^2 + \|X\Delta_{H_0}\|^2$$

$$\begin{array}{c} \sim \\ \sim X_{n-p}^2 \\ \sim X_r^2 \end{array}$$

which leads to

$$F = \frac{(SSE_0 - SSE_1) / r}{SSE_1 / (n-p)} \sim F_{r, n-p}$$

$\sim \dots$  = previous Wald-test-statistic

### Alternative approach via substitution

Partitioning  $C$  and  $\beta$  by columns and rows we can write

$$C\beta = d$$

as

$$\left[ \begin{array}{cc} C_0 & C_r \\ \sim & \sim \\ r \times (p-r) & r \times r \end{array} \right] \left[ \begin{array}{c} \beta_0 \\ \beta_r \end{array} \right] = d$$

"redundant"

where  $C_r$  is invertible.

Using blockwise matrix multiplication

$$C_0 \beta_0 + C_r \beta_r = d$$

such that

$$\beta_r = C_r^{-1} (C_0 \beta_0 - d)$$

Partitioning  $X$  the same way, we can write

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{\varepsilon}$$

as

$$\underline{Y} = \begin{bmatrix} X_0 & X_r \end{bmatrix} \begin{bmatrix} \underline{\beta}_0 \\ \underline{\beta}_r \end{bmatrix} + \underline{\varepsilon}$$

$n \times (p+r)$     $n \times r$

$$= X_0 \underline{\beta}_0 + X_r \underline{\beta}_r + \underline{\varepsilon}$$

$$= X_0 \underline{\beta}_0 + X_r C_r (\underline{C}_0 \underline{\beta}_0 - \underline{d}) + \underline{\varepsilon}$$

$$= (X_0 + X_r C_r \underline{C}_0) \underline{\beta}_0 + X_r C_r \underline{d} + \underline{\varepsilon}$$

$\underbrace{(X_0 + X_r C_r \underline{C}_0)}$   
modified design  
matrix  $X^*$

$\underbrace{X_r C_r \underline{d}}$   
offset  
 $\underline{a}$

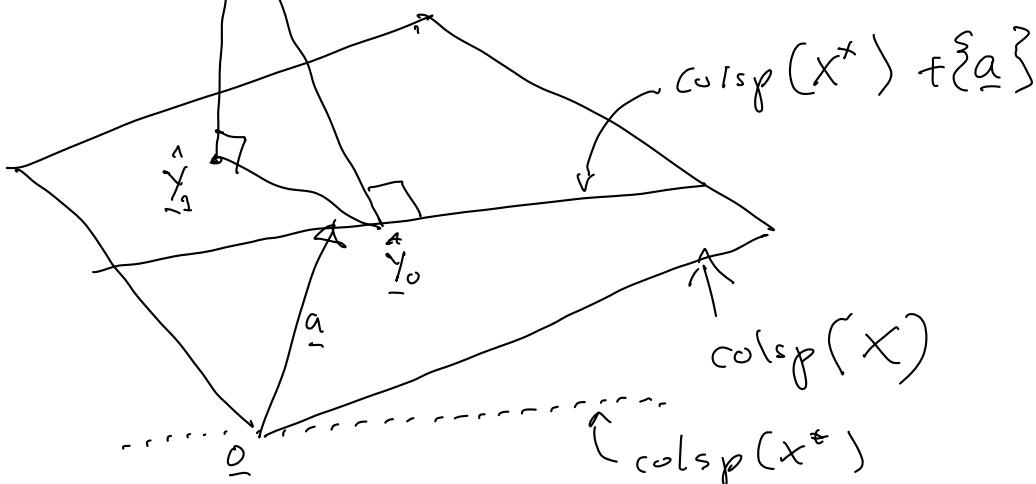
2.9

Fitted values  $\hat{Y}_0$  under  $H_0$  given by

projection of  $\underline{Y}$  onto affine subspace

$$\{X\underline{\beta} \mid \underline{\beta} \in \mathbb{R}^n, C\underline{\beta} = \underline{d}\} = \text{colsp}(X^*) + \{\underline{a}\}$$

Minkowski sum



7.9

Thus

$$\hat{y}_0 = \text{Proj}_{\text{colsp}(X^*)}(\underline{y} - \underline{a}) + \underline{a}$$

$$\hat{\beta}_0 = (X^{*\top} X^*)^{-1} X^{*\top} (\underline{y} - \underline{a})$$

and

$$\hat{\beta}_r = -C_r^{-1} (C_0 \hat{\beta}_0 - \underline{d})$$

and the same F-statistic as before.

This can be generalized to fitting GLMs under linear equality constraints facilitating LRTs (as opposed to less reliable Wald-tests).

## Categorical covariates (factors) and interactions

One-way anova,  $k$  groups: For  $j=1, 2, \dots, k$ ,  
and  $i=1, 2, \dots, n_j$

$$Y_{ij} = \mu + \underbrace{\alpha_j}_{\text{effect of belonging to group } j} + \varepsilon_{ij}, \quad \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$$

Alternative notation (see Wood, p. 198).

For  $i=1, 2, \dots, n$  ( $n = \sum_{j=1}^k n_j$ )

$$Y_i = \mu + \alpha_{j(i)} + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

Unknown parameters:  $\mu, \alpha_1, \alpha_2, \dots, \alpha_k, \sigma^2$

In matrix notation

$$Y = \begin{matrix} j=1 \\ j=2 \\ j=3 \end{matrix} \left[ \begin{matrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{matrix} \right] \begin{matrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \beta \end{matrix} + \varepsilon$$

Columns are not indep. (model is not identifiable)

Omit second column  $\Leftrightarrow$  imposing the constraint  $\alpha_1 = 0$ .  
(treatment contrasts)

(Symbolic (Wilkinson-Rogers) notation in R :)

$$\gamma \sim \underbrace{\text{glucose}}_{\text{factor}}$$

Thus, for

$$j=1, E(y_i) = \mu$$

$$j=2, E(y_i) = \mu + \alpha_2$$

$$j=3, E(y_i) = \mu + \alpha_3$$

i.e.  $\alpha_2$  is the difference in  $E(y_i)$  in group 2 relative to group 1 (the control/placebo treatment)

### Interactions (with numerical covariate)

"Alternative" notation: For  $j=1, 2, \dots, n$

$$Y_i = \mu + \underbrace{\alpha_{j(i)}}_{\substack{\text{Main effect} \\ \text{of belonging} \\ \text{to group } j}} + \underbrace{\beta x_i}_{\substack{\text{Main effect} \\ \text{of num.} \\ \text{covariate } x_i}} + \underbrace{\gamma_{j(i)} \cdot x_i}_{\substack{\text{Interaction} \\ \text{between} \\ j(i) \text{ and } x_i}} + \varepsilon_i$$

$$= (\underbrace{\mu + \alpha_{j(i)}}_{\substack{\text{"Intercept"} \\ \text{in group } j}}) + (\underbrace{\beta + \gamma_{j(i)}}_{\substack{\text{"Slope"} \text{ in} \\ \text{group } j}}) x_i + \varepsilon_i$$

Symbolic notation

$$\log(\text{dian}) \sim \text{glucose} + \log(\text{cons}) + \text{glucose} : \log(\text{cons})$$

Impose the constraints  $\alpha_1 = 0$  and  $\beta_1 = 0$

(if using treatment contrasts), the model in matrix notation becomes

$$y_i = \begin{cases} j=1 \\ j=2 \\ j=3 \end{cases} \left[ \begin{array}{cccccc} 1 & 0 & 0 & x_1 & 0 & 0 \\ 1 & 0 & 0 & x_2 & 0 & 0 \\ 1 & 1 & 0 & x_3 & x_3 & 0 \\ 1 & 1 & 0 & x_4 & x_4 & 0 \\ 1 & 0 & 1 & x_5 & 0 & x_5 \\ 1 & 0 & 1 & x_6 & 0 & x_6 \end{array} \right] \begin{bmatrix} \mu \\ \alpha_2 \\ \alpha_3 \\ \beta \\ \beta_2 \\ \beta_3 \end{bmatrix} + \varepsilon_i$$

Diagram illustrating the structure of the matrix equation:

- The matrix has 6 columns: 1, 0, 0,  $x_1$ , 0, 0; 1, 0, 0,  $x_2$ , 0, 0; 1, 1, 0,  $x_3$ ,  $x_3$ , 0; 1, 1, 0,  $x_4$ ,  $x_4$ , 0; 1, 0, 1,  $x_5$ , 0,  $x_5$ ; 1, 0, 1,  $x_6$ , 0,  $x_6$ .
- Braces on the left indicate groups:  $j=1$  covers the first two columns;  $j=2$  covers the third and fourth columns;  $j=3$  covers the fifth and sixth columns.
- Braces on the right indicate groups:  $\mu$ ,  $\alpha_2$ ,  $\alpha_3$ ,  $\beta$ ,  $\beta_2$ ,  $\beta_3$ .
- A bracket at the bottom indicates the error term  $\varepsilon_i$ .
- A large bracket at the bottom indicates the overall structure of the equation.
- A handwritten 'X' is drawn under the last three columns of the matrix.

Thus, for

$$j=1, E(y_i) = \mu + \beta x_1$$

$$j=2, E(y_i) = \mu + \alpha_2 + (\beta + \beta_2)x_1$$

$$j=3, E(y_i) = \mu + \alpha_j + (\beta + \delta_j)x_i$$

↑                      ↓

difference in  
intercept and slope  
relative to the  
control group  $j=1$

Sum to zero contrasts: Impose the constraints

$$\sum_{j=1}^k \alpha_j = 0, \quad \sum_{j=1}^k \gamma_j = 0$$

such that e.g.

$$\alpha_3 = -\alpha_1 - \alpha_2, \quad \gamma_3 = -\gamma_1 - \gamma_2$$

R: contrasts(glucose)  $\leftarrow$  "contr. sum"

$\underbrace{\phantom{000}}_{\text{factor}}$

Model in matrix notation:

$$\underline{y} = \begin{cases} j=1 \\ j=2 \\ j=3 \end{cases} \left[ \begin{array}{cccccc} 1 & 1 & 0 & x_1 & x_1 & 0 \\ 1 & 1 & 0 & x_2 & x_2 & 0 \\ 0 & 1 & 1 & x_3 & 0 & x_3 \\ 0 & 0 & 1 & x_4 & 0 & x_4 \\ 1 & -1 & -1 & x_5 & -x_5 & -x_5 \\ 1 & -1 & -1 & x_6 & -x_6 & -x_6 \end{array} \right] \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \beta \\ \delta_1 \\ \delta_2 \end{bmatrix} + \varepsilon$$

$\underbrace{\quad\quad\quad}_{X}$

Second alternative notation:

$$\log(\text{diam}_i) = \mu + \alpha_{\text{glucose}(i)} + \beta \log(\text{conc}_i) + \gamma_{\text{glucose}(i)} \cdot \log(\text{conc}_i) + \varepsilon_i$$

GLMs: Linear predictor constructed the same way.

# Binary regression (Ch. 5)

Sept. 1

Ungrouped data:  $(y_i, x_i)$ ,  $i=1, 2, \dots, n$

Model:

$$y_i \sim \text{Bernoulli}(\pi_i)$$

where  $\pi_i$  depends on a linear predictor

$y_i = x_i^\top \beta$  via response function

$$\pi_i = h(y_i)$$

$$g(\pi_i) = y_i$$

Logit link

$$g(\pi) = \ln \frac{\pi}{1-\pi}, \quad h(y) = \frac{e^y}{1+e^y} = \frac{1}{1+e^{-y}}$$

Odds interpretation

$$\frac{\pi_i}{1-\pi_i} = e^{\beta_0 + \beta_1 x_{i1} + \cdots + \beta_k x_{ik}}$$

Unit change in  $x_{i1}$  changes the odds

by a factor (odds ratio) of  $e^{\beta_1}$

## Probit link

Response function

$$\pi = \Phi(\gamma) = \int_{-\infty}^y \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz = \underbrace{\text{pnorm}(ct)}_{R\text{-funktn}}$$

Link function

$$\Phi^{-1}(\pi) = \gamma$$

$$\text{probit}(\pi) = \gamma = \text{qnorm}(p^i)$$

Example : Model :

Age at time of menarche

$$T_i \sim N(\mu, \sigma^2)$$

We then have

$$\pi_i = P(Y_i = 1) = P(T_i < t_i) \quad \begin{matrix} \underbrace{} \\ \text{age at medical} \\ \text{examination} \end{matrix}$$

$\sim N(0, 1)$

$$= P\left(\frac{T_i - \mu}{\sigma} < \frac{t_i - \mu}{\sigma}\right)$$

$$= \Phi\left(\frac{t_i - \mu}{\sigma}\right)$$

$$\Phi^{-1}(\pi_i) : \text{probit } \pi_i = \frac{t_i - \mu}{\sigma}$$

a glm.

$$= -\underbrace{\mu}_{\beta_0} + \underbrace{\frac{1}{\sigma} t_i}_{\beta_1}$$

glm( $y \sim \text{age}$ , family = binomial(link = "probit"))

It follows that

$$\sigma = \frac{1}{\beta_1} \quad \mu = -\frac{\beta_0}{\beta_1}$$

Hence, the MLEs of  $\mu$  and  $\sigma$  are

$$\hat{\sigma} = \frac{1}{\hat{\beta}_1} = \frac{1}{0.86} = 1.16 \text{ years}, \quad \hat{\mu} = -\frac{\hat{\beta}_0}{\hat{\beta}_1} = -\frac{-11.32}{0.86} =$$

$$= 13.22 \text{ years}$$

Exercise 5. Compute the S.E.s.  
using the delta method

$$Z = f(x, Y) \approx f(\mu_x, \mu_Y) + \frac{\partial f}{\partial x}(x - \mu_x) + \frac{\partial f}{\partial Y}(Y - \mu_Y)$$

$$\text{Var} Z = \left( \frac{\partial f}{\partial x} \right)^2 \text{Var} X + \left( \frac{\partial f}{\partial y} \right)^2 \text{Var} Y + 2 \frac{\partial f}{\partial x} \frac{\partial f}{\partial y} \text{Cov}(X, Y)$$

7.9

9.9

### Complementary log-log link

Response func.

$$-e^{-y}$$

$$\pi = h(y) = 1 - e^{-y}$$

and link func.

$$g(\pi) = \underbrace{\ln(-\ln(1-\pi))}_{\text{cloglog}} = y$$

Project 1 : Show that this corresponds to the latent variable model

$$T_i \sim \text{Weibull}(\alpha, \beta) \text{ and } \pi_i = P(T_i > t_i)$$

Example 1: Survival analysis:  
Suppose that the hazard for the  $i$ -th subject is

$$\lambda_i(t) = \lambda_0(t) e^{\underbrace{\beta_1 x_i}_{\text{baseline hazard}}}$$

$$= \lim_{\Delta t \rightarrow 0} \frac{P(T_i \leq t + \Delta t \mid T_i \geq t)}{\Delta t}$$

At the end of the experiment at time  $T$   
each subject is dead with probability

$$\pi_i = P(Y_i = 1) = P(T_i < T)$$

$$= 1 - e^{- \int_0^T \lambda_i(t) dt} \quad \left( \text{known from survival analysis} \right)$$

$$= 1 - e^{- e^{\beta_0} \int_0^T \lambda_0(t) dt} \quad \underbrace{e^{\beta_0}}$$

$$= 1 - e^{- e^{\beta_0 + \beta_1 x_i}}$$

$$= 1 - e^{- e^{\eta_i}}$$

i.e.

$$\text{log}(\pi_i) = \beta_0 + \beta_1 x_i$$

Interpretation of  $\beta_1$ : One unit change in  $x_i$   
changes the hazard  $\lambda_i(t)$  by a factor  
of  $e^{\beta_1}$ , e.g. if  $\beta_1 = -0.1$  the  $e^{\beta_1} \approx 0.91$ ,  
that is, a 9% decrease.

Example 2: Presence-absence data in ecology:

Suppose the number of individuals of a given species  $Z_i \sim \text{Poisson}(\lambda_i)$  at given sampling sites  $i = 1, 2, \dots, n$ , but that we only observe/record the binary response

$$Y_i = \begin{cases} 1 & \text{if } Z_i \geq 1 \text{ (presence)} \\ 0 & \text{if } Z_i = 0 \text{ (absent)} \end{cases}$$

Assume that  $\ln \lambda_i = \gamma_i = \underline{x}_i^\top \beta$

Then

$$\pi_i = P(Y_i = 1) = P(Z_i \geq 1)$$

$$= 1 - P(Z_i = 0)$$

$$= 1 - \frac{\lambda_i^0 e^{-\lambda_i}}{0!}$$

$$= 1 - e^{-e^{\gamma_i}}$$

that is,

$$\text{cloglog}(\pi_i) = \gamma_i = \underline{x}_i^\top \beta$$

Grouped data : Each observation  $i = 1, 2, \dots, n$  involves  $n_i$  Bernoulli trials such that

$$Y_i \sim \text{bin}(n_i, \pi_i)$$

Notation :

$$\bar{Y}_i = \frac{Y_i}{n_i}$$

$$\bar{Y}_i \sim \underbrace{\text{bin}(n_i, \pi_i)}_{n_i}$$

R :  $\text{glm}(\text{cbind}(y, n-y) \sim x, \text{binomial})$

$\underbrace{\quad}_{n \times 2 \text{ matrix}}$

$$\left[ \begin{array}{cc} Y_1 & n_1 - Y_1 \\ Y_2 & n_2 - Y_2 \\ \vdots & \vdots \end{array} \right] \quad \left\{ \begin{array}{l} n \text{ rows} \end{array} \right.$$

If  $\text{Var}(Y_i) > n_i \pi_i (1 - \pi_i)$  we have overdispersion.

# Log-likelihood (non-grouped data)

$$l(\beta) = \sum_{i=1}^n l_i(\beta)$$

$$P(Y_i = y_i) = \pi_i^{y_i} (1 - \pi_i)^{1 - y_i}$$

where the contributions

$$l_i(\beta) = y_i \ln \pi_i + (1 - y_i) \ln (1 - \pi_i)$$

$$= y_i \ln \frac{\pi_i}{1 - \pi_i} + \ln (1 - \pi_i)$$

canonical

link func.

$$= y_i \gamma_i + \ln \frac{1}{1 + e^{\gamma_i}}$$

$$\pi_i = \frac{e^{\gamma_i}}{1 + e^{\gamma_i}}$$

$$1 - \pi_i = \frac{1}{1 + e^{\gamma_i}}$$

$$= y_i \gamma_i - \ln (1 + e^{\gamma_i})$$

$$\gamma_i = \hat{x}_i^\top \beta$$

Score function

$$s(\beta) = \frac{\partial}{\partial \beta} l(\beta) = \frac{\partial}{\partial \beta} \sum_{i=1}^n l_i(\beta) = \sum_{i=1}^n \underbrace{\frac{\partial}{\partial \beta} l_i(\beta)}_{s_i(\beta)}$$

Contribution from  $i$ -th observation

$$s_i(\beta) = \frac{\partial}{\partial \beta} l_i(\beta)$$

$$= \frac{\partial}{\partial \beta} (y_i \gamma_i - \ln (1 + e^{\gamma_i}))$$

$$= y_i \frac{\partial \gamma_i}{\partial \beta} - \frac{e^{\gamma_i}}{1 + e^{\gamma_i}} \cdot \frac{\partial \gamma_i}{\partial \beta}$$

$$= \left( y_i - \pi_i \right) \frac{\partial}{\partial \beta} \left( x_i^\top \beta \right)$$

$$= (y_i - \pi_i) x_i$$

$$= (y_i - \pi_i) \begin{bmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \end{bmatrix}$$

Total score function

$$\underbrace{s(\beta)}_{p \times 1} = \sum_{i=1}^n s_i(\beta) = \sum_{i=1}^n (y_i - \pi_i) \underbrace{x_i}_{p \times 1}$$

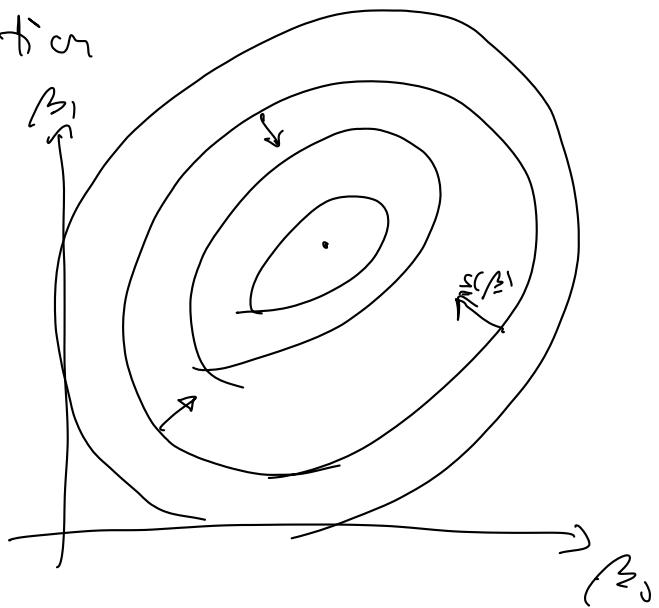
Keep in mind that  $\pi_i$  is a function of  $y_i$  which in turn is a function of  $\beta$ , i.e.  $\pi_i = \pi_i(y_i(\beta))$ .

Need to solve the equation

$$\underbrace{s(\beta)}_{p \times 1} = 0$$

A system of  $p$  non-linear eqs. in  $p$  unknowns.

Next time: algorithm for solving this numerically



## Newton's method

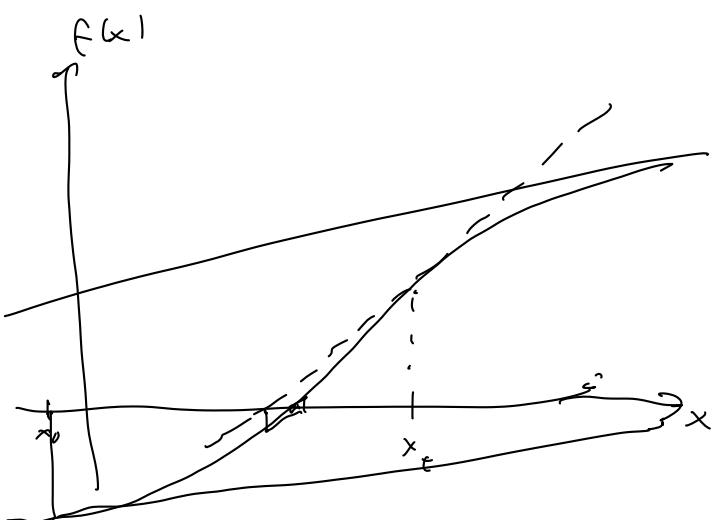
In one dimension. Find the root of

$$f(x) = 0$$

Idea: Current best guess at the solution  $x_t$ .

Set linear approximation around  $x_t$  equal to zero and solve for  $x$ :

$$f(x) \approx f(x_t) + f'(x_t)(x - x_t) = 0$$



$$x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)}$$

Multiple dimensions: Current best guess is  $\beta_t$

Approximate  $\underline{s}(\beta)$  by

$$\underline{s}(\beta) \approx \underline{s}(\beta_t) + \frac{\partial}{\partial \beta} \underline{s}^T(\beta_t) (\beta - \beta_t)$$

$$= \underline{s}(\beta_t) + \frac{\partial}{\partial \beta} \frac{\partial}{\partial \beta} \underline{l}(\beta_t) (\beta - \beta_t)$$

$$= \underline{s}(\beta_t) - H(\beta_t) (\beta - \beta_t) = 0$$

Solving for  $\beta$  we obtain

$$\underline{\beta}_{t+1} = \underline{\beta}_t + \left( H^{-1}(\underline{\beta}_t) \right)^{-1} S(\underline{\beta}_t) \quad (\text{Newton's method})$$

In practice we replace  $H(\underline{\beta}_t)$  by  $F(\underline{\beta}_t)$  (expected Fisher information)

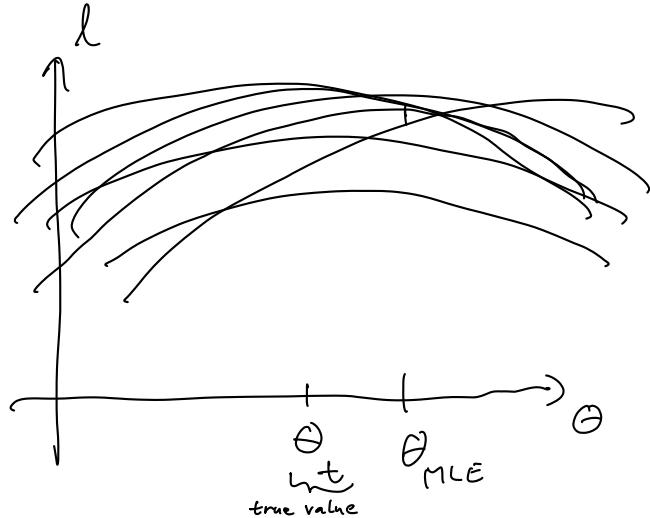
$$\underline{\beta}_{t+1} = \underline{\beta}_t + F^{-1}(\underline{\beta}_t) S(\underline{\beta}_t) \quad \left. \begin{array}{l} \text{Fisher scoring} \\ \text{algorithm} \end{array} \right\}$$

Recall that

$$F(\underline{\beta}) = E(H(\underline{\beta})) = E\left(-\frac{\partial^2 l}{\partial \underline{\beta} \partial \underline{\beta}^T}\right)$$

Note:  
 $l(\underline{\beta}; \mathbf{y})$   
is a random  
variable

Some general properties of the log-likelihood



$$\begin{aligned} l(\theta; \mathbf{Y}) &\stackrel{\text{random}}{=} \ln L[\theta; \mathbf{Y}] \\ &= \ln f(\mathbf{Y}; \theta) \\ &\underbrace{\qquad\qquad\qquad}_{\text{joint pdf of all data } \mathbf{Y}} \end{aligned}$$

$$1. \quad E\left(\frac{\partial l}{\partial \underline{\theta}} \Big|_{\theta=\theta_t}\right) = E(S(\underline{\theta}_t)) = 0$$

Proof:

$$E\left(\frac{\partial l}{\partial \underline{\theta}}\right) = \int \frac{\partial}{\partial \underline{\theta}} \ln f(y) f(y) dy$$

$$\begin{aligned}
 &= \int \frac{1}{f(y)} \cdot \frac{\partial f}{\partial \underline{\theta}} \cdot \cancel{f(y)} dy \\
 &= \frac{\partial}{\partial \underline{\theta}} \int f(y) dy \\
 &\Rightarrow \frac{\partial}{\partial \underline{\theta}} 1 = 0
 \end{aligned}$$

by def.

$$\begin{aligned}
 2. \text{Var} \left( \underbrace{\frac{\partial l}{\partial \underline{\theta}}}_{p \times 1} \right) &= E \left( \frac{\partial l}{\partial \underline{\theta}} - E \left( \frac{\partial l}{\partial \underline{\theta}} \right) \right) \left( \frac{\partial l}{\partial \underline{\theta}} - E \left( \frac{\partial l}{\partial \underline{\theta}} \right) \right)^T \\
 &\stackrel{p \times p}{=} E \left( \frac{\partial l}{\partial \underline{\theta}} \frac{\partial l}{\partial \underline{\theta}^T} \right)_{p \times p}
 \end{aligned}$$

$$3. \text{Var}\left(\left.\frac{\partial l}{\partial \underline{\theta}}\right|_{\underline{\theta}=\underline{\theta}_t}\right) = E\left(\left.-\frac{\partial^2 l}{\partial \underline{\theta} \partial \underline{\theta}^\top}\right|_{\underline{\theta}=\underline{\theta}_t}\right) = F(\underline{\theta}_t)$$

$\underbrace{\phantom{\dots}}_{P \times P}$        $\underbrace{\phantom{\dots}}_{P \times P}$

Proof:

$$\frac{\partial}{\partial \underline{\theta}} E\left(\frac{\partial l}{\partial \underline{\theta}^\top}\right) = \underbrace{\frac{\partial}{\partial \underline{\theta}}}_{\substack{P \times 1 \\ 1 \times P}} \underbrace{\underline{\theta}^\top}_{P \times P} = 0$$

$$\frac{\partial}{\partial \underline{\theta}} \int \frac{\partial}{\partial \underline{\theta}^\top} \ln f(y) \cdot f(y) dy = 0$$

$$\int \frac{\partial}{\partial \underline{\theta}} \left( \frac{\partial}{\partial \underline{\theta}^\top} \ln f(y) \cdot f(y) \right) dy = 0$$

$$\begin{aligned} \int & \left( \frac{\partial}{\partial \underline{\theta}} \frac{\partial}{\partial \underline{\theta}^\top} \ln f(y) \cdot f(y) + \underbrace{\frac{\partial}{\partial \underline{\theta}} f(y)}_{= \frac{\partial}{\partial \underline{\theta}} \ln f(y) \cdot f(y)} \frac{\partial}{\partial \underline{\theta}^\top} \ln f(y) \cdot f(y) \right) dy = 0 \\ & - F(\underline{\theta}) + E(s(\underline{\theta}) s^\top(\underline{\theta})) = 0 \end{aligned}$$

$$F(\underline{\theta}) = \text{Var}(s(\underline{\theta}))$$

## Fisher information for binary stem

Recall:

$$s(\beta) = \sum_{i=1}^n s_i(\beta) = \sum_{i=1}^n (y_{i\cdot} - \hat{u}_{i\cdot}) \underline{x}_{i\cdot}$$

$$\hat{u}_{i\cdot} = \hat{u}_i(\eta_i(\beta))$$

Expected Fisher information via property 3:

$$F(\beta) = \text{Var}(s(\beta)) = \text{Var}\left(\sum_{i=1}^n (y_{i\cdot} - \hat{u}_{i\cdot}) \underline{x}_{i\cdot}\right)$$

indep obs.

$$= \sum_{i=1}^n \text{Var}\left(\underbrace{\underline{x}_{i\cdot} \cdot y_{i\cdot}}_{p \times 1}\right)$$

$$= \sum_{i=1}^n \underline{x}_{i\cdot} \underbrace{\text{Var}(y_{i\cdot})}_{\hat{u}_{i\cdot}(1-\hat{u}_{i\cdot})} \underline{x}_{i\cdot}^\top$$

$$= \sum_{i=1}^n \hat{u}_{i\cdot}(1-\hat{u}_{i\cdot}) \underbrace{\underline{x}_{i\cdot} \underline{x}_{i\cdot}^\top}_{p \times p} \underbrace{\hat{u}_{i\cdot}(1-\hat{u}_{i\cdot})}_{p \times 1}$$

$$p \times p$$

Observed Fisher information:

$$H_i(\beta) = - \frac{\partial^2}{\partial \beta \partial \beta^T} \ell_i(\beta_i) = - \frac{\partial^2}{\partial \beta^2} S_i^T(\beta)$$

$$= - \frac{\partial}{\partial \beta} \left( y_i - \frac{e^{y_i}}{1 + e^{y_i}} \right) \cdot \underline{x}_i^T$$

$$\text{II} \quad \frac{e^{y_i} (1 + e^{y_i}) - e^{y_i} e^{y_i}}{(1 + e^{y_i})^2} \cdot \frac{\partial y_i}{\partial \beta} \cdot \underline{x}_i^T$$

$$= \pi_i (1 - \pi_i) \frac{\partial}{\partial \beta} \left( \underline{x}_i^T \beta \right) \cdot \underline{x}_i^T$$

$$= \pi_i (1 - \pi_i) \underline{x}_i \underline{x}_i^T \cdot \frac{e^{y_i}}{(1 + e^{y_i})^2}$$

$$= \underbrace{\frac{1}{1 + e^{y_i}}}_{1 - \pi_i} \cdot \underbrace{\frac{e^{y_i}}{1 + e^{y_i}}}_{\pi_i}$$

$$H(\beta) = \sum_{i=1}^n \pi_i (1 - \pi_i) \underline{x}_i \underline{x}_i^T$$

Expected Fisher info. via def.:

$$F(\beta) = E H(\beta) = E \left( \sum_{i=1}^n \pi_i^* (1 - \pi_i^*) \underline{x}_i \underline{x}_i^T \right)$$

and so

$$F(\beta) = H(\beta)$$

Because we are using  
the canonical link  
function.

9.9  
14.9

Is  $F(\beta)$  always positive definite and invertible?  
 ( p. 283 in Fahrmeir )

Binary regression:

$$F(\beta) = \sum_{i=1}^n \pi_i(1-\pi_i) \underline{x}_i \underline{x}_i^\top = \underline{X}^\top W \underline{X}$$

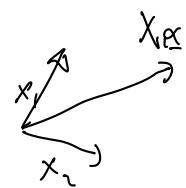
where  $W = \text{diag}(\pi_i(1-\pi_i))$

$$= \begin{bmatrix} \pi_1(1-\pi_1) & & \\ & \pi_2(1-\pi_2) & \\ & & \ddots \end{bmatrix}$$

Recall: If  $\underline{X}$  has full rank ( $p \leq n$  indep columns), then

$$\|\underline{X}\underline{\alpha}\| = \underline{\alpha}^\top \underline{X}^\top \underline{X}\underline{\alpha} > 0 \quad \text{for all } \underline{\alpha} \neq \underline{0}$$

$\underbrace{\text{lin. comb.}}_{\text{of columns}} \text{ of } \underline{X}$



i.e.  $\underline{X}^\top \underline{X}$  is positive definite and invertible.

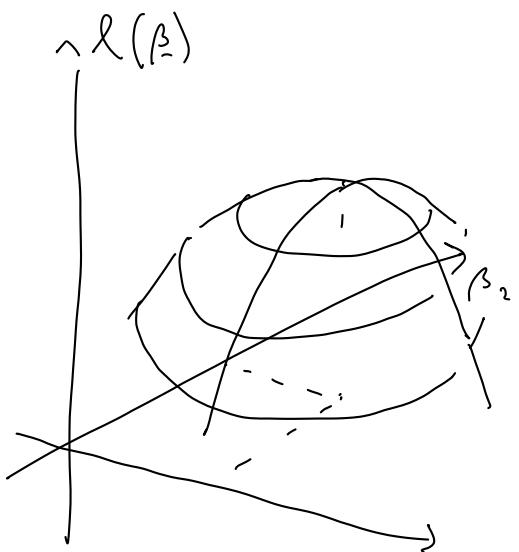
Let  $Z = \text{diag}(\sqrt{\pi_i(1-\pi_i)}) \underline{X}$ , i.e. each row  $i=1, \dots, n$  of  $\underline{X}$  multiplied by  $\sqrt{\pi_i(1-\pi_i)}$ . For  $\beta \in \mathbb{R}^p$ ,

$0 < \pi_i < 1$  and  $\sqrt{\pi_i(1-\pi_i)} > 0$  for all  $i=1, 2, \dots, n$

and  $\text{rowsp}(Z) = \text{rowsp}(\underline{X}) \Rightarrow \text{rank}(Z) = \text{rank}(\underline{X})$ ,

$\Rightarrow F(\beta) = Z^\top Z$  is positive definite and invertible.

If  $H(\beta) = F(\beta) \Rightarrow l(\beta)$  is concave (with a single maxima)



But  $F(\beta)$  may tend to positive semi-definite as some  $\beta_i \rightarrow +/- \infty$

Recommended exercise: For a binary glm with an intercept term included and a logit link-function, show that

$$\frac{1}{n} \sum y_i = \frac{1}{n} \sum \hat{\pi}_i \Leftrightarrow \sum_{i=1}^n (y_i - \hat{\pi}_i) = 0$$

where  $\hat{\pi}_i = h(\hat{y}_i) := \frac{1}{1 + e^{-x_i^\top \beta}}$  (the fitted probabilities).

Would the same be true for other link functions / models without an intercept?

## Fisher scoring algorithm for the linear model

$$y_i \sim N(\mu_i, \sigma^2) \quad \underbrace{\mu_i = \eta_i}_{\text{identity link}}, \quad \eta_i = \underline{x}_i^\top \underline{\beta}$$

$$\ell_i = C - \frac{1}{2\sigma^2} (y_i - \mu_i)^2$$

$$\begin{aligned} s_i &= \frac{\partial \ell_i}{\partial \beta} = \frac{1}{\sigma^2} (y_i - \mu_i) \frac{\partial \mu_i}{\partial \beta} \\ &= \frac{1}{\sigma^2} (y_i - \mu_i) \underline{x}_i \end{aligned}$$

$$s(\underline{\beta}) = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \mu_i) \underline{x}_i = \dots = \frac{1}{\sigma^2} \underline{x}^\top (Y - X\underline{\beta})$$

$$F(\underline{\beta}) = \text{Var}(s(\underline{\beta})) = \left(\frac{1}{\sigma^2}\right)^2 \sum_{i=1}^n \text{Var}(\underline{x}_i | Y)$$

$$= \frac{1}{\sigma^4} \sum x_i \underbrace{\text{Var}(y_i)}_{\text{Var}} \underline{x}_i^\top$$

$$= \frac{1}{\sigma^2} \sum_{i=1}^n x_i^T x_i$$

$$= \frac{1}{\sigma^2} X^T X$$

Fish. scor. algorithm.

$$\beta_{t+1} = \beta_t + \left( \frac{X^T X}{\sigma^2} \right)^{-1} \underbrace{\left( X^T Y - X^T X \beta_t \right)}_{S(\beta_t)}$$

$F(\beta_t)$

$$= \beta_t + (X^T X)^{-1} X^T Y - \beta_t$$

$$= (X^T X)^{-1} X^T Y = \hat{\beta}_{MLE}$$

Very fast convergence (since  $\ell(\beta)$  is exactly quadratic in  $\beta$ !)

Fisher scoring alg. sometimes diverge

Minimal example: A single observation, no covariates  
(intercept only)

$$y \sim \text{bin}(n, \pi), \log(\pi) = \gamma, \gamma = \beta_0$$

$$\ell(\beta) = \ln \binom{n}{y} + y \ln \pi + (n-y) \ln(1-\pi)$$

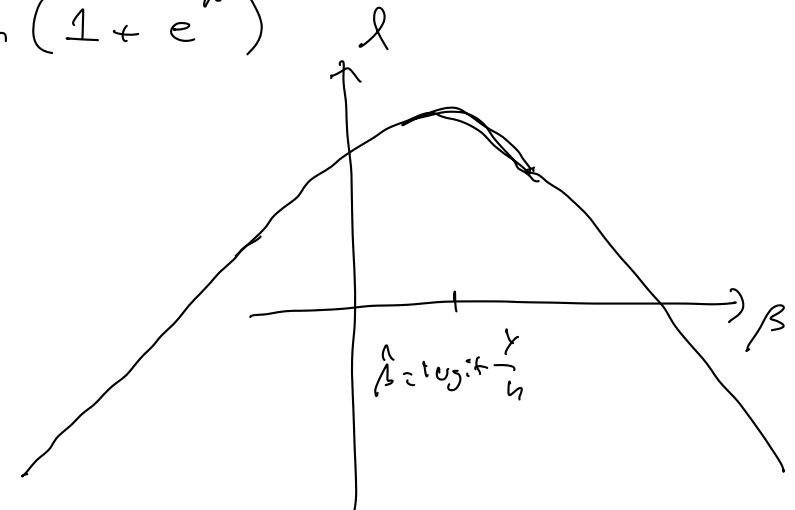
$$= C + y \ln \frac{\pi}{1-\pi} + n \ln(1-\pi)$$

$$= C + y\beta + n \ln \left( 1 - \frac{e^\beta}{1+e^\beta} \right) = \dots =$$

$$= C + y\beta - n \ln(1+e^\beta)$$

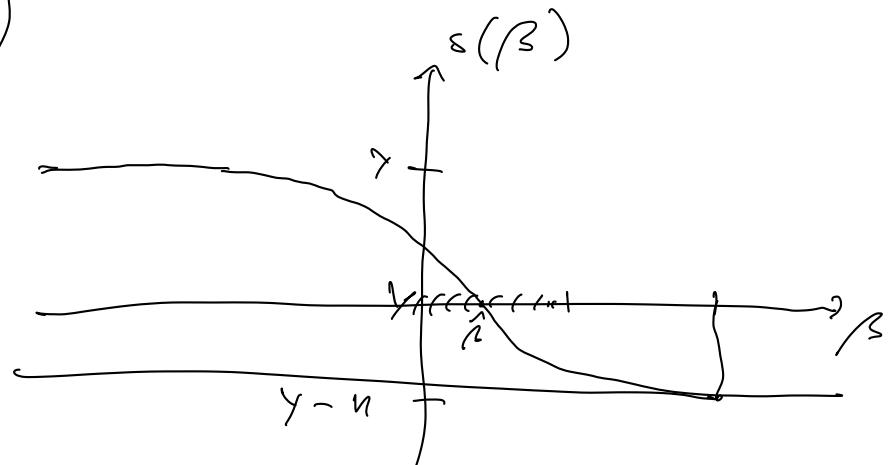
$$\hat{\pi} = \frac{y}{n}$$

$$\hat{\beta} = \log \frac{y}{n}$$



$$s(\beta) = y - n \frac{e^\beta}{1+e^\beta}$$

Newton method  
= Fisher scoring alg.  
in this case  
will overshoot for bad  
starting values.



Starting values sometimes important, `glm( , start= )`

## Some asymptotic properties of MLEs

With independent observations, as the sample size  $n \rightarrow \infty$ ,

$$\hat{\theta}_{\text{MLE}} \sim N(\underline{\theta}, F^{-1}(\underline{\theta}))$$

Sketch of proof: Linear approximation of  $s(\hat{\theta})$  around  $\underline{\theta}$  (the true value of  $\theta$ ):

$$\underbrace{s(\hat{\theta})}_{=0} \approx s(\underline{\theta}) + \left. \frac{\partial^2 l}{\partial \underline{\theta} \partial \underline{\theta}^T} \right|_{\underline{\theta}} \cdot (\hat{\theta} - \underline{\theta})$$

$$= H(\underline{\theta})$$

Hence

$$\hat{\theta} - \underline{\theta} \approx H^{-1}(\underline{\theta}) s(\underline{\theta})$$

and

$$\sqrt{n}(\hat{\theta} - \underline{\theta}) \approx \sqrt{n} H^{-1}(\underline{\theta}) s(\underline{\theta})$$

$$= \left( \frac{H(\underline{\theta})}{n} \right)^{-1} \frac{s(\underline{\theta})}{\sqrt{n}} \xrightarrow{n \rightarrow \infty} \left( \frac{F(\underline{\theta})}{n} \right)^{-1} w$$

where  $w \sim N(0, \frac{F(\underline{\theta})}{n})$  (since  $E(s(\underline{\theta})) = 0$  and  $\text{Var}(s(\underline{\theta})) = F(\underline{\theta})$ )

Hence,

$$E(\hat{\theta}) \approx \underline{\theta}$$

$$\text{Var}(\sqrt{n}(\hat{\theta} - \underline{\theta})) \approx \left( \frac{F(\underline{\theta})}{n} \right)^{-1} \frac{F(\underline{\theta})}{n} \left( \frac{F(\underline{\theta})}{n} \right)$$

$$= n F^{-1}(\underline{\theta})$$

$$\text{Var}(\hat{\theta}) \approx F^{-1}(\underline{\theta})$$

In practice (for finite sample size) we can use this as an approximation and estimate  $\text{Var}(\hat{\theta})$  by

$$\text{Var}(\hat{\theta}) = \hat{F}'(\theta) = F^{-1}(\hat{\theta})$$

or by  $H'(\hat{\theta})$  (as suggested by Efron & Hinkley, 1978)

summary(glm(...))  
vcov(...)

## Hypothesis testing for GLMs

$$H_0: C\beta = d \quad \text{vs.} \quad H_1: C\beta \neq d$$

$\underbrace{\phantom{0}}_{r \times p}$     $\underbrace{\phantom{0}}_{p \times 1}$

| Fitted model         | Type of test   |
|----------------------|--|
| Both $H_0$ and $H_1$ | <p>Likelihood ratio test (usually most accurate)</p> <p>Under <math>H_0</math>:</p> $2 \left( \ell(\hat{\beta}_1) - \ell(\hat{\beta}_0) \right) \stackrel{\text{asymp.}}{\sim} \chi^2_r \quad r = p_1 - p_0$ <p>R: anova(mod0, mod1, test = "chisq")<br/>or drop1(mod1, test = "chisq")</p>  |
| Only $H_1$           | <p>Wald test: Under <math>H_0</math>,</p> $(C\hat{\beta}_1 - d)^T (C F'(\hat{\beta}_1) C^T)^{-1} (C\hat{\beta}_1 - d) \stackrel{\text{asymp.}}{\sim} \chi^2_r$ <p>For the case <math>H_0: \beta_i = 0</math> vs. <math>H_1: \beta_i \neq 0</math> the Wald test simplifies to <math>\frac{\hat{\beta}_i^2}{\text{Var}(\hat{\beta}_i)} \stackrel{\text{asymp.}}{\sim} \chi^2_1</math> or <math>Z = \frac{\hat{\beta}_i}{\sqrt{\text{Var}(\hat{\beta}_i)}} \stackrel{\text{asymp.}}{\sim} N(0, 1)</math></p> |

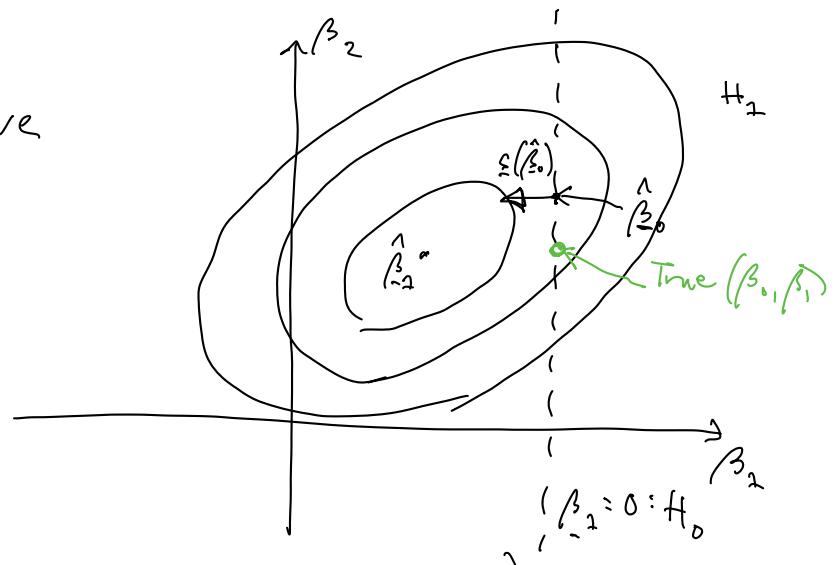
summarizing (mod 1) gives you  
Wald - tests or "Z - tests"

Only  $H_0$

Score test (under  $H_0$ , P - r parameters)

$$\underbrace{S_1(\hat{\beta}_0)}_{1 \times p}^T \underbrace{F_1^{-1}(\hat{\beta}_0)}_{p \times p} \underbrace{S_1(\hat{\beta}_0)}_{p \times 2} \stackrel{\text{asympt}}{\sim} \chi_r^2$$

$S_1(\hat{\beta}_0)$  will have  
r nonzero  
elements



anova (mod 0, test = "Rao")

# Deviance (and testing goodness-of-fit)

16.9

Plays the same role as SSE of LMs.

Def.: The deviance  $D$  of a fitted model is as

$$D = 2 \left( l_{\text{saturated}} - l(\hat{\beta}) \right)$$

where  $l(\hat{\beta})$  is the maximum log likelihood under the fitted model ( $p_0$  parameters) and  $l_{\text{saturated}}$  is the maximum log likelihood under a model with  $P_1 = n$  parameters ( $n$  is no. of observations)

If the candidate fitted model is true, then

$$D \stackrel{\text{asympt}}{\sim} \chi^2_{n - p_0}$$

$$E(D) = n - p_0$$

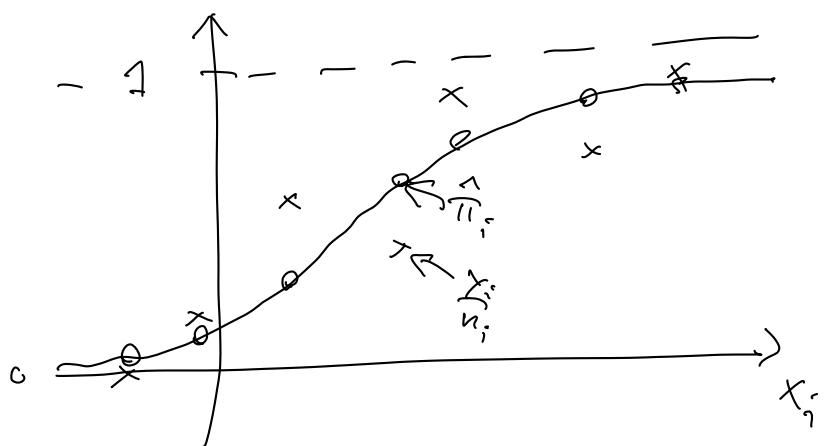
Reject  $H_0$  if  $D > \chi^2_{n-p_0, \alpha}$  (Goodness-of-fit test)

Deviance for binary regression model (grouped data)

$$Y_i \sim \text{bin}(n_i, \pi_i)$$

$$D = 2 \left[ l\left(\frac{y_1}{n_1}, \frac{y_2}{n_2}, \dots\right) - l\left(\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_k\right) \right]$$

$$\hat{\pi}_i = h(x_i^T \beta)$$



$$= 2 \sum_{i=1}^n y_i \ln \frac{y_i}{n_i} + (n_i - y_i) \ln \left(1 - \frac{y_i}{n_i}\right) - y_i \ln \hat{\pi}_i - (n_i - y_i) \ln (1 - \hat{\pi}_i)$$

$$= 2 \sum_{i=1}^n y_i \ln \frac{\bar{y}_i}{\hat{\pi}_i} + (n_i - y_i) \ln \frac{1 - \bar{y}_i}{1 - \hat{\pi}_i} \quad \bar{y}_i = \frac{y_i}{n_i}$$

$$= \sum_{i=1}^n n_i^2$$

deviance  
residuals

## Two types of residuals (binary regression)

Pearson residuals

$$r_i = \frac{y_i - n_i \hat{\pi}_i}{\sqrt{n_i \hat{\pi}_i (1 - \hat{\pi}_i)}}$$

Alternative goodness-of-fit statistic

$$\sum r_i^2 \stackrel{\text{approx}}{\sim} \chi_{n-p}^2$$

Deviance residuals:

$$r_i = \underbrace{\text{sign}(\bar{y}_i - \hat{\pi}_i)}_{\text{sign}(x)} \sqrt{2 \left( y_i \ln \frac{\bar{y}_i}{\hat{\pi}_i} + (n_i - y_i) \ln \frac{1 - \bar{y}_i}{1 - \hat{\pi}_i} \right)}$$

$$\text{sign}(x) = \begin{cases} 1 & \text{for } x > 0 \\ 0 & \text{for } x = 0 \\ -1 & \text{for } x < 0 \end{cases}$$

$$E(r_i) \approx 0, \quad \text{Var}(r_i) \approx 1$$

Estimating the dispersion parameter (in quasi-likelihood models)

$$\hat{\phi} = \frac{D}{n-p}$$

or

$$\hat{\phi} = \frac{1}{n-p} \sum_{\substack{i=1 \\ \text{pearson residual}}} r_i^2$$

$\left\{ \begin{array}{l} \text{used by} \\ \text{quasibinomial /} \\ \text{quasipoisson} \end{array} \right\}$

# Asymptotic distribution of the LRT - statistic (Wood sec. 4.4)

Want to test

$$H_0: \underbrace{R(\underline{\theta}) = 0}_{r \times 1} \quad \left. \begin{array}{l} r \text{ nonlinear} \\ \text{equality constraints} \end{array} \right.$$

vs.

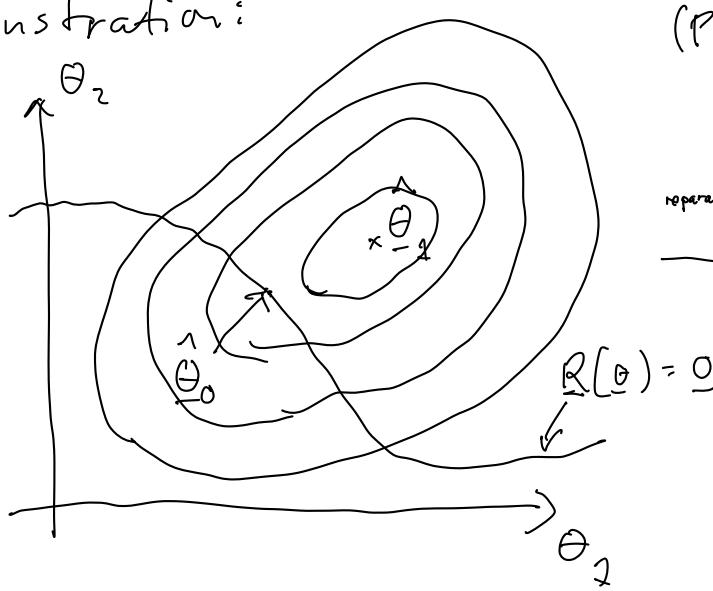
$$H_1: R(\underline{\theta}) \neq 0$$

Theorem: Under  $H_0$ , asympt.

$$\text{LRT} = 2[\ell(\hat{\underline{\theta}}_1) - \ell(\hat{\underline{\theta}}_0)] \sim \chi_r^2$$

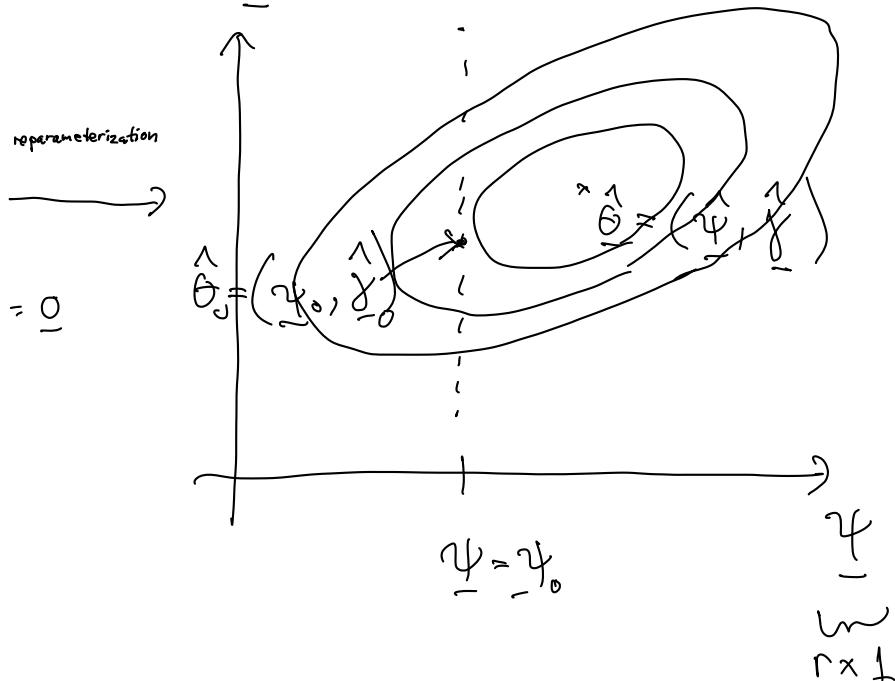
as  $n \rightarrow \infty$  where  $\hat{\underline{\theta}}_1$  and  $\hat{\underline{\theta}}_0$  are the MLE of  $\underline{\theta}$  under  $H_1$  and  $H_0$ .

Illustration:



$$(p-r \times 1) \quad \underline{\gamma}$$

reparameterization



Prof: Asymptotic approximation of  $\ell(\underline{\theta})$  around the MLE  $\hat{\underline{\theta}} = (\hat{\psi}, \hat{\gamma})$  under  $H_2$ .

$$\ell(\underline{\theta}) \approx \ell(\hat{\underline{\theta}}) - \frac{1}{2} (\underline{\theta} - \hat{\underline{\theta}})^T H (\underline{\theta} - \hat{\underline{\theta}}), \quad (1)$$

$$-\frac{1}{2} (\underline{\theta} - \hat{\underline{\theta}})^T H (\underline{\theta} - \hat{\underline{\theta}})$$

$$L(\underline{\theta}) \approx L(\hat{\underline{\theta}}) e$$

$$\propto MVN\left(\hat{\underline{\theta}}, \begin{bmatrix} \Sigma_{\psi\psi} & \Sigma_{\psi\gamma} \\ \Sigma_{\gamma\psi} & \Sigma_{\gamma\gamma} \end{bmatrix}\right)$$

$\underbrace{\qquad\qquad\qquad}_{= H^{-1}}$

$$H = -\frac{\partial^2 \ell}{\partial \underline{\theta} \partial \underline{\theta}^T}$$

Hence, under  $H_0$ ,  $L(\underline{\theta})$  is maximised by  $\hat{\underline{\theta}}_0 = (\hat{\psi}_0, \hat{\gamma}_0)$

where  $\hat{\gamma}_0 = \hat{\gamma} + \sum_{\psi\gamma} \Sigma_{\psi\psi}^{-1} (\hat{\psi}_0 - \hat{\psi}) \quad (2)$

Analogous to the formula  $E(\underline{x}_1 | \underline{x}_2) = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (\underline{x}_2 - \mu_2)$  when  $\begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \end{bmatrix} \sim MVN$ .

Since  $\Sigma = H^{-1}$ , partitions of these matrices satisfies

$$\begin{bmatrix} \Sigma_{\psi\psi} & \Sigma_{\psi\gamma} \\ \Sigma_{\gamma\psi} & \Sigma_{\gamma\gamma} \end{bmatrix} \begin{bmatrix} H_{\psi\psi} & H_{\psi\gamma} \\ H_{\gamma\psi} & H_{\gamma\gamma} \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & I \end{bmatrix}$$

which leads to 3 eqs. in  $\Sigma_{\psi\psi}, \Sigma_{\gamma\psi} = \Sigma_{\psi\gamma}^T, \Sigma_{\gamma\gamma}$ :

$$\sum_{\psi\psi} H_{\psi\psi} + \sum_{\psi\gamma} H_{\psi\gamma} = I \quad \sum_{\psi\psi} H_{\psi\gamma} + \sum_{\gamma\gamma} H_{\gamma\gamma} = 0$$

$$\sum_{\psi\psi} H_{\psi\psi} - \sum_{\psi\gamma} H_{\psi\gamma} H_{\gamma\gamma}^{-1} = P \quad \sum_{\psi\gamma} = - \sum_{\psi\psi} H_{\psi\gamma} H_{\gamma\gamma}^{-1}$$

with solution

$$H_{\psi\psi} - H_{\psi\gamma} H_{\gamma\gamma}^{-1} H_{\gamma\psi} = \sum_{\psi\psi} \quad (4) \quad \sum_{\gamma\psi} = - H_{\gamma\gamma}^{-1} H_{\gamma\psi} \sum_{\psi\psi} \quad (3)$$

Substitute (3) into (2)

$$\hat{\gamma}_0 = \hat{\gamma} - H_{\gamma\gamma}^{-1} H_{\gamma\psi} \sum_{\psi\psi} (\psi_0 - \hat{\psi})$$

$$\hat{\gamma}_0 - \hat{\gamma} = - H_{\gamma\gamma}^{-1} H_{\gamma\psi} (\psi_0 - \hat{\psi})$$

Using (1), the log likelihood under  $H_0$  is then

$$\ell(\psi_0, \hat{\gamma}_0) \approx \ell(\hat{\psi}, \hat{\gamma}) - \underbrace{\frac{1}{2} \left[ \begin{bmatrix} \psi_0 - \hat{\psi} \\ \hat{\gamma}_0 - \hat{\gamma} \end{bmatrix}^T H \begin{bmatrix} \psi_0 - \hat{\psi} \\ \hat{\gamma}_0 - \hat{\gamma} \end{bmatrix} \right]}_{= LRT}$$

and

$$LRT \approx (\psi_0 - \hat{\psi})^T \begin{bmatrix} I & \\ -H_{\gamma\gamma}^{-1} H_{\gamma\psi} & \end{bmatrix}^T \begin{bmatrix} H_{\psi\psi} & H_{\psi\gamma} \\ H_{\gamma\psi} & H_{\gamma\gamma} \end{bmatrix} \begin{bmatrix} I & \\ -H_{\gamma\gamma}^{-1} H_{\gamma\psi} & \end{bmatrix} (\psi_0 - \hat{\psi})$$

$$= \dots = (\underline{\psi}_0 - \hat{\underline{\psi}})^\top \left( H_{44}^{-1} - H_{48} H_{88}^{-1} H_{84} \right) (\underline{\psi}_0 - \hat{\underline{\psi}})$$

(4)

$$= (\underline{\psi}_0 - \hat{\underline{\psi}})^\top \sum_{44}^{-1} (\underline{\psi}_0 - \hat{\underline{\psi}})$$

As  $n \rightarrow \infty$ ,  $\text{Var}(\hat{\underline{\theta}}) \simeq F^{-1}(\underline{\theta}) \simeq H^{-1}(\underline{\theta}) = \sum$

and in particular

$$\text{Var}(\hat{\underline{\psi}}) \simeq \sum_{44}$$

Hence, under  $H_0$ ,

$$\text{LRT} = 2 \left( \ell(\hat{\underline{\theta}}_1) - \ell(\hat{\underline{\theta}}_0) \right) \simeq (\underline{\psi}_0 - \hat{\underline{\psi}})^\top \sum_{44}^{-1} (\underline{\psi}_0 - \hat{\underline{\psi}}) \sim \chi_n^2$$

# Numerical issues in computation of the deviance (project 1, prob. 1b)

Log likelihood:

$$l(\hat{\beta}) = \sum_{i=1}^n y_i \ln \hat{\lambda}_i + (1-y_i) \ln (1-\hat{\lambda}_i)$$

Ok for computing  $l(\hat{\beta})$  and  $f_0$ , but  
under the saturated model

Contribution to the log likelihood

$$l_i(\hat{\lambda}_i) = \ln f(y_i; \hat{\lambda}_i)$$

$$\approx \ln \frac{\hat{\lambda}_i^{y_i} e^{-\hat{\lambda}_i}}{y_i!}$$

$$= y_i \ln \hat{\lambda}_i - \hat{\lambda}_i - \ln y_i!$$

Under the saturated model  $\hat{\lambda}_i = y_i$ . For observing

$$l_i(\hat{\lambda}_i) = 0 \ln 0 - 0 - \ln 0!$$

$$= 0 \cdot \text{NaN} - 0 - \ln 1$$

$$= \text{NaN}$$

But

$$\ln f(0; 0) = \ln 1 = 0$$

dpois(0, 0, log = TRUE)

## Nested models

Tests based on change in deviance

$$D_0 - D_1 = \dots = LRT \sim \chi^2_{p_1 - p_0}$$

why?

## Models selection

- Based on hypothesis testing

Include only significant covariates

Not recommended

- Choose the model with the smallest

Akaike information criteria

$$AIC = -2l(\hat{\theta}) + 2p$$

↴              ↴  
 measure      penalize  
 model fit      complex  
 models

This is  
an estimate of Kullback - Leibler distance  
to the true model (next week)

16.19

21.9

Example:

library (I8UR)  
data (eba 1977) } See R-code on course web page  
...

$Y_i$  (cases) Number of lung cancer cases

$n_i$  (pop) sub pop. size

$j(i)$  (age) age categories (6 levels)

$k(i)$  (city) Fredericia, Horsens, Kolding, Vejle

Aim: Is there an elevated risk of lung cancer in Fredericia caused by petrochemical factory?

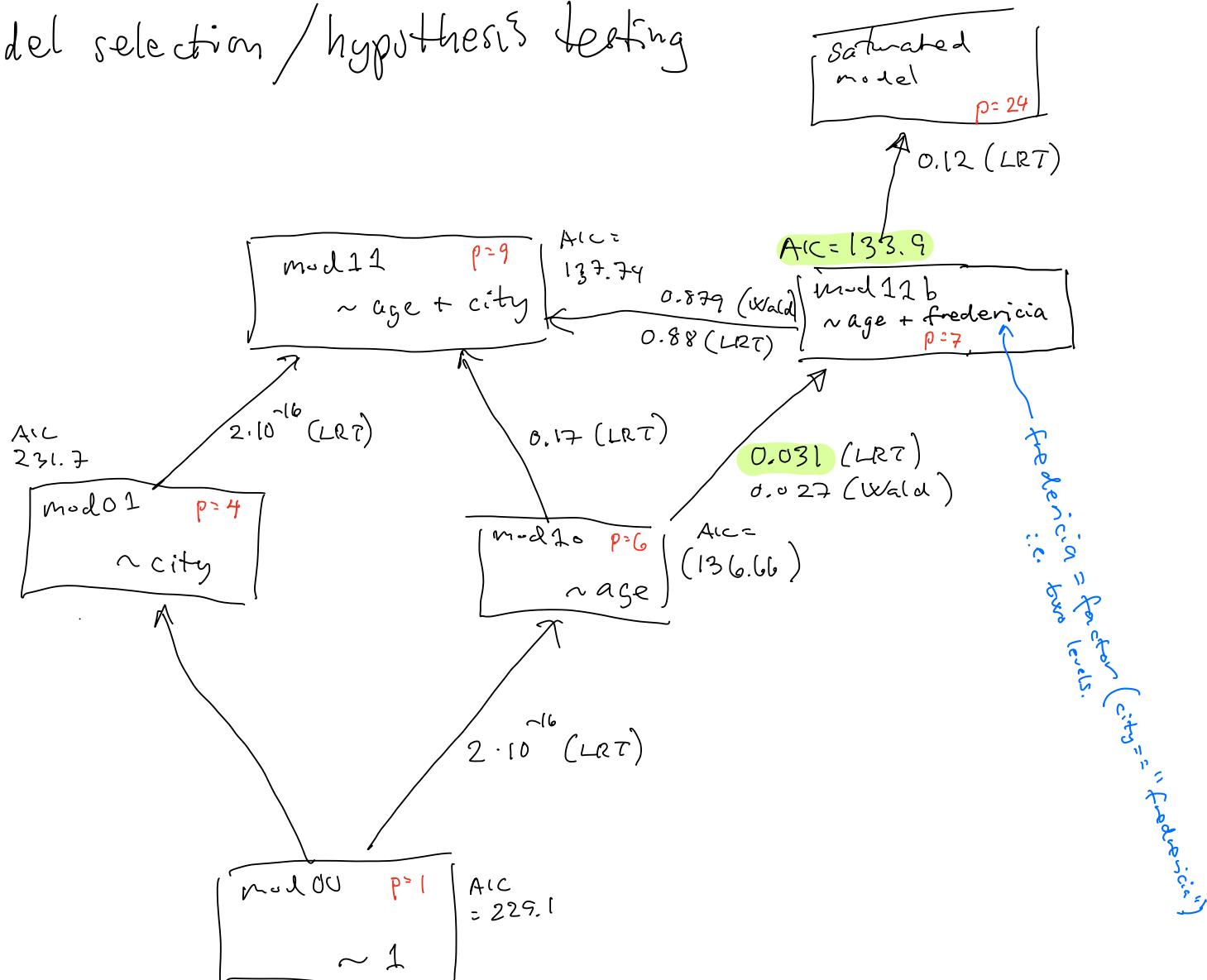
Model (modell):

$$Y_i \sim \text{bin}(n_i, \pi_i)$$

$$\text{logit } \pi_i = \mu + \underbrace{\alpha_{j(i)}}_{\text{effect of age}} + \underbrace{\beta_{k(i)}}_{\text{effect of city}}$$

$$\beta = (\mu, \alpha_1, \alpha_2, \alpha_3, \dots, \alpha_6, \beta_2, \beta_3, \beta_4)^T \quad (p=9 \text{ parameters})$$

# Model selection / hypothesis testing



Want to fit the model  $\beta_1 \neq \beta_2 = \beta_3 = \beta_4$  (mod11b)

Testing this via a Wald test

$$C\beta = d$$

$$\begin{bmatrix} 0 & \dots & 0 & 0 & 1 & -1 & 0 \\ 0 & \dots & 0 & 0 & 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} \cdot \\ \vdots \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Wald test:  $(C\hat{\beta} - d)^T (C F^{-1}(\hat{\beta}) C^T)^{-1} (C\hat{\beta} - d) \sim \chi^2_2$

from mod11

## Theory behind AIC (Wood sec. 4.6)

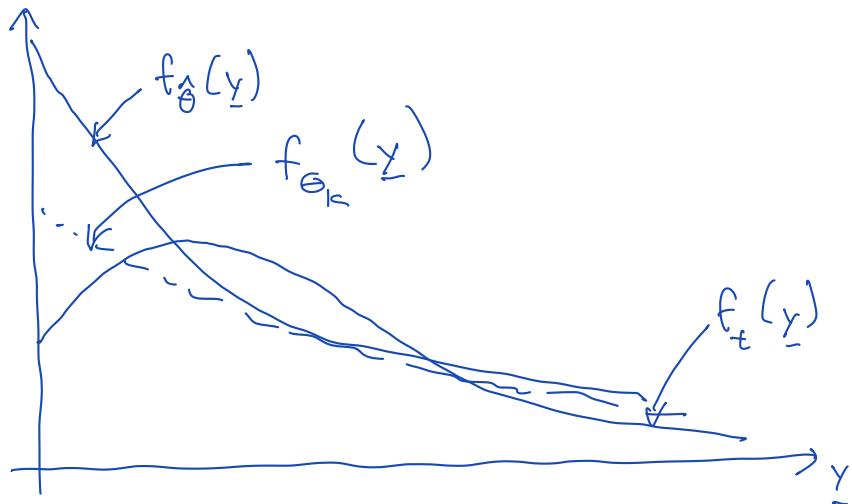
Aim: Choose a candidate model  $f_\theta(\underline{y})$  minimising the Kullback - Leibler distance to the true model  $f_t(\underline{y})$

$$K(f_\theta, f_t) = \int (\ln f_t(\underline{y}) - \ln f_\theta(\underline{y})) f_t(\underline{y}) d\underline{y} \quad (1)$$

$$= E(\ln f_t(\underline{Y}) - \ln f_\theta(\underline{Y})) , \quad \underline{Y} \sim f_t$$

One-dimensional illustrations:

" $\underline{Y}$  distributed as a random variable with pdf (or pmf)  $f_t$ "



Problems:  $f_t$  is unknown and  $\theta$  is only estimated by  $\hat{\theta} = \hat{\theta}(\underline{Y}_2)$ ;  $\underline{Y}_2 \sim f_t$  and indep. of  $\underline{Y}$ . However,

$$E K(f_{\hat{\theta}}, f_t) = E E(\ln f_t(\underline{Y}) - \ln f_{\hat{\theta}}(\underline{Y}) \mid \underline{Y}_2)$$

can be estimated up to a constant only involving  $f_t$  as follows:

First consider  $K(f_\theta, f_t)$ : This is minimized by some  $\theta = \theta_K$  satisfying

$$\frac{\partial}{\partial \theta} \int_{\theta > \theta_k} (\ln f_t(y) - \ln F_\theta(y) f_t(y) dy) = 0$$

$$\int_{\theta = \theta_k} \left. \frac{\partial}{\partial \theta} \ln F_\theta(y) \right|_{\theta = \theta_k} f_t(y) dy = 0 \quad (2)$$

Approximating  $\ln f_\theta(y)$  around  $\theta_k$ ,

$$\ln f_\theta(y) \approx \ln F_{\theta_k}(y) + (\hat{\theta} - \theta_k)^\top \left. \frac{\partial}{\partial \theta} \ln f_\theta(y) \right|_{\theta = \theta_k} + \frac{1}{2} (\hat{\theta} - \theta_k)^\top \left. \frac{\partial^2 \ln f_\theta(y)}{\partial \theta \partial \theta^\top} \right|_{\theta = \theta_k}$$

Substituting this into (1) we then find that

$$\begin{aligned} K(f_{\hat{\theta}(y)}, f_t) &\approx \left( \left( \ln f_t(y) - \ln F_{\theta_k}(y) - (\hat{\theta} - \theta_k)^\top \left. \frac{\partial}{\partial \theta} \ln f_\theta(y) \right|_{\theta = \theta_k} \right) + \frac{1}{2} (\hat{\theta} - \theta_k)^\top \left. \frac{\partial^2 \ln f_\theta(y)}{\partial \theta \partial \theta^\top} \right|_{\theta = \theta_k} (\hat{\theta} - \theta_k)^\top f_t(y) dy \right) \\ &\stackrel{(1), (2)}{=} K(f_{\theta_k}, f_t) + \underbrace{\frac{1}{2} (\hat{\theta} - \theta_k)^\top F_{\theta_k}^\top (\hat{\theta} - \theta_k)}_{(*)} \end{aligned}$$

For  $f_\theta$  sufficiently close to  $f_t$   $\xrightarrow{(*)}$  (3)

$$\hat{\theta} \stackrel{\text{approx}}{\sim} N(\theta_k, F_{\theta_k}^{-1})$$

such that  $(*) \sim \chi_p^2$  and

$$E K(f_{\hat{\theta}}, f_t) \approx K(f_{\theta_k}, f_t) + \frac{1}{2} p. \quad (4)$$

Next consider

$$\begin{aligned}
 E(-\ell(\hat{\theta})) &= E\left(-\ell(\theta_k) - (\ell(\hat{\theta}) - \ell(\theta_k))\right) \\
 &\stackrel{(3)}{\approx} E(-\ell(\theta_k)) - \frac{1}{2} p \\
 &= \int (\ln f_t(y) - \ln f_{\theta_k}(y)) f_t(y) dy - \frac{1}{2} p \\
 &\quad - \int \ln f_t(y) f_t(y) dy \\
 &= K(f_{\theta_k}, f_t) - \frac{1}{2} p - \int \ln f_t(y) f_t(y) dy \tag{5}
 \end{aligned}$$

Eliminating  $K(f_{\theta_k}, f_t)$  from (4) and (5) we obtain

$$EK(f_{\hat{\theta}}, f_t) \approx E(-\ell(\hat{\theta})) + p + \text{constant term only involving } f_t(y)$$

Thus, an estimate of  $EK(f_{\hat{\theta}}, f_t)$  (up to an unknown constant) is

$$EK(f_{\hat{\theta}}, f_t) \approx -\ell(\hat{\theta}) + p,$$

and an estimate of twice of this is

$$AIC = -2\ell(\hat{\theta}) + 2p.$$

Note the similarity between  $K(f_{\hat{\theta}}, f_t)$  and the leave-out-one log-likelihood cross-validation score

$$CV = -\sum_{i=1}^n \ln f_{\hat{\theta}_{-i}}(y_i)$$

# Poisson regression with identity non-canonical link function

Model

$$y_i \sim \text{Poisson}(\lambda_i), \quad \underbrace{\lambda_i = \eta_i}_{\text{identity link}}, \quad \eta_i = \underline{x}_i^\top \underline{\beta}$$

Likelihood

$$L(\underline{\beta}) = \prod_{i=1}^n \frac{\lambda_i^{y_i} e^{-\lambda_i}}{y_i!}$$

Log-lik.

$$\begin{aligned} \ell(\underline{\beta}) &= \sum_{i=1}^n y_i \ln \lambda_i - \lambda_i - \ln y_i! \\ &= \sum_{i=1}^n y_i \ln \eta_i - \eta_i - \ln y_i! \end{aligned}$$

Score function

$$\begin{aligned} s(\underline{\beta}) &= \sum_{i=1}^n y_i \frac{1}{\eta_i} \frac{\partial \eta_i}{\partial \underline{\beta}} - \frac{\partial \eta_i}{\partial \underline{\beta}} \\ &= \sum_{i=1}^n \left( \frac{y_i}{\eta_i} - 1 \right) \underline{x}_i. \end{aligned}$$

Recall that  
 $\frac{\partial \eta_i}{\partial \underline{\beta}} = \frac{\partial}{\partial \underline{\beta}} (\underline{x}_i^\top \underline{\beta}) = \underline{x}_i$

Fisher info

$$F(\underline{\beta}) = \text{Var}(s(\underline{\beta})) = \sum_{i=1}^n \underline{x}_i \frac{\text{Var}(y_i)}{y_i^2} \underline{x}_i^\top$$

$$\begin{aligned} &= \sum_{i=1}^n \underline{x}_i \frac{\lambda_i}{y_i^2} \underline{x}_i^\top = \underline{X}^\top \text{diag}\left(\frac{1}{\lambda}\right) \underline{X} \\ &\quad \underbrace{\underline{x} \quad \underline{x}^\top \quad \lambda}_{p \times 1 \quad 1 \times p \quad 1 \times p} \end{aligned}$$

$$\text{diag}\left(\frac{1}{\lambda}\right) = \begin{bmatrix} \frac{1}{\lambda_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{\lambda_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \cdots & \frac{1}{\lambda_n} \end{bmatrix}$$

Observed Fisher info:

$$\begin{aligned}
 H(\beta) &= -\frac{\partial^2}{\partial \beta \partial \beta^\top} l(\beta) = -\underbrace{\frac{\partial}{\partial \beta} \left( \frac{\partial}{\partial \beta^\top} l(\beta) \right)}_{p \times 1 \quad 1 \times p} = -\frac{\partial}{\partial \beta} \delta^\top(\beta) \\
 &= -\sum \frac{\partial}{\partial \beta_i} \left( \frac{y_i}{\eta_i} - 1 \right) x_i^\top \\
 &= +\sum \frac{y_i}{\eta_i^2} \frac{\partial \eta_i}{\partial \beta_i} x_i^\top \\
 &\leq +\sum \underbrace{\frac{y_i}{\eta_i^2} x_i x_i^\top}_{p \times p} = X^\top \text{diag}\left(\frac{y}{\eta^2}\right) X = Z^\top Z \\
 &\neq F(\beta)
 \end{aligned}$$

Even if  $X$  has full rank  $p < n$ ,

$H(\beta)$  is possibly only semi-positive definite. Why?

Let

$$Z = \text{diag}\left(\sqrt{\frac{y}{\eta^2}}\right) X \quad \left( \begin{array}{l} \text{rows of } X \text{ gets multiplied} \\ \text{by } \sqrt{y_i/\eta_i^2} = 0 \text{ for some obs.} \end{array} \right)$$

$Z$  may not have  $p$  independent rows (and columns)  
and  $\text{rank}(Z) \leq p$ . Thus  $H(\beta) = Z^\top Z$  may only be  
positive semi-definite and non-invertible since

$$(Z\alpha)^\top (Z\alpha) = \underbrace{\|Z\alpha\|^2}_{\text{any linear comb. of columns in } Z} \geq 0$$

any linear comb. of columns in  $Z$ .

One reason for preferring Fisher scoring alg. over Newton's method.

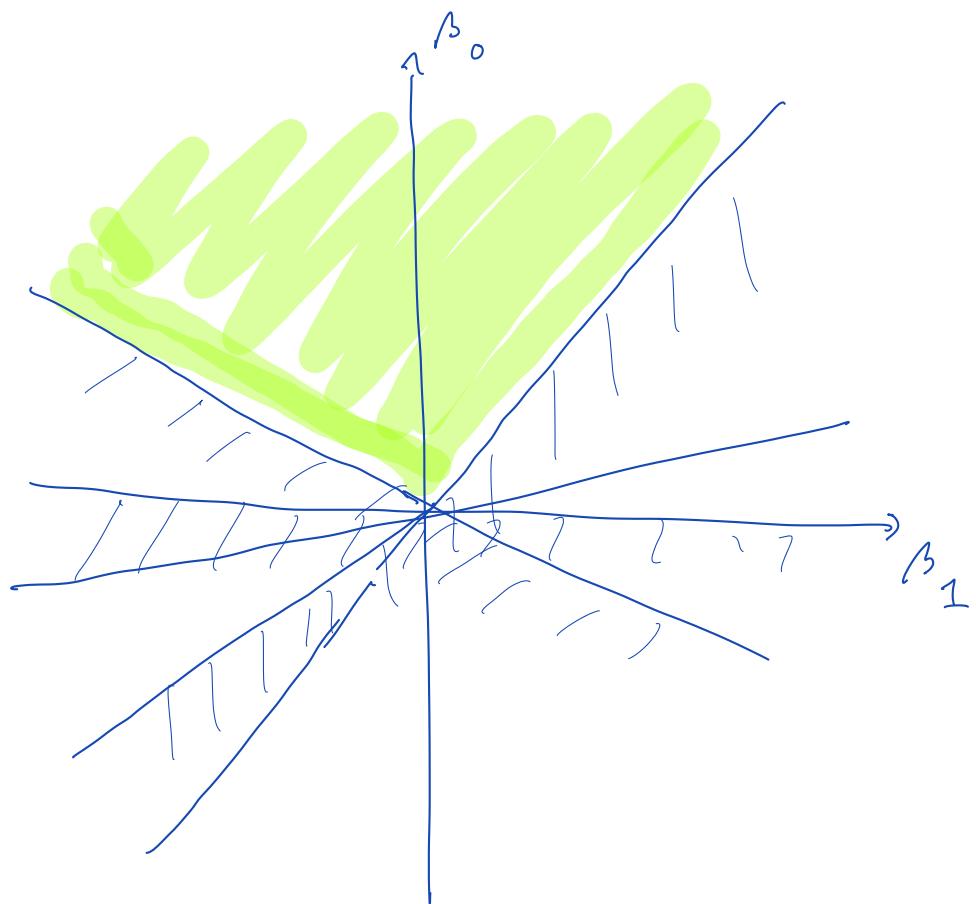
23.9

Constraints on  $\beta$ : For  $i = 1, 2, \dots, n$  and  $p=2$

$$\lambda_i = \gamma_i = \beta_0 + \beta_1 x_i \geq 0$$

$$\beta_0 \geq -\beta_1 x_i$$

$$\beta_0 \geq \max_i (-\beta_1 x_i) = \begin{cases} \beta_1 \max_i (-x_i) & \text{for } \beta_1 > 0 \\ -\beta_1 \max_i (x_i) & \text{for } \beta_1 \leq 0 \end{cases}$$



## GLMs for positive continuous responses (Ch. 5.3)

Gamma

$$f(y|a,b) = \frac{b^a}{\Gamma(a)} y^{a-1} e^{-by}, \quad EY = \frac{a}{b}, \quad \text{Var } Y = \frac{a}{b^2}$$

Accommodates right skew.

Belongs to exponential family since

$$\ln f(y|a,b) = a \ln b + (a-1) \ln y - by - \ln \Gamma(a)$$

$$\begin{aligned} &= -bx + a \ln b \\ &\quad \cancel{a} \quad \cancel{b}(\theta) \\ &= \underbrace{-\frac{b}{a}y + \ln(b/a)}_{1/\phi} - \underbrace{\frac{\ln(1/a)}{1/a}}_{1/\phi} - \ln \Gamma(a) + (a-1) \ln y \\ &\quad \quad \quad c(y, \phi, w) \end{aligned}$$

For the parameterization  $\mu = \frac{a}{b}$  and  $v = a$  then

$$\theta = -\frac{1}{\mu} \quad \text{and} \quad b(\theta) = -\ln(-\theta) \quad \phi = \frac{1}{a} = \frac{1}{v}$$

GLM with canonical link function

$$y_i \sim \text{Gamma}(\mu_i, v)$$

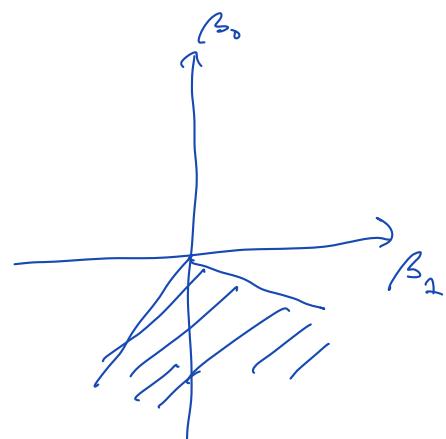
where

$$-\frac{1}{\mu_i} = \eta_i = \bar{x}_i^\top \beta$$

But

$$\mu_i > 0$$

$$-\frac{1}{\mu_i} = \eta_i = \bar{x}_i^\top \beta < 0$$



If  $\eta_i = \beta_0 + \beta_1 x_i$ , then  $\beta_0 < \min_i (-\beta_1 x_i)$

Gamma GLM with non-canonical log-link

$$\log(\mu_i) = \eta_i = \hat{x}_i^T \hat{\beta}$$

$$Y_i \sim \text{Gamma}(\mu_i, \nu)$$

ensures that

$$\mu_i = e^{\eta_i} = e^{\hat{x}_i^T \hat{\beta}} \geq 0 \quad \text{for all } \hat{\beta} \in \mathbb{R}^P$$

Note that  $\ln Y_i$  becomes strongly left skewed  
(exercise 10)

Alternative lognormal GLM with a log-link

$$y_i \sim \text{lognormal}(\mu_i, \sigma)$$

$$\log(\mu_i) = \hat{x}_i^T \hat{\beta}$$



$$\ln y_i = \hat{x}_i^T \hat{\beta} + \varepsilon_i \quad (\text{a LM})$$

## 5.4. General treatment of GLMs

Log likelihood contribution on the exponential family form:

$$l_i(\beta) = \ln f(x_i | \theta_i, \phi, w_i)$$

$$= \frac{w_i}{\phi} \left( y_i \theta_i - b(\theta_i) \right) + c(y_i, w_i, \phi).$$

Recall that

$$\mu_i = b'(\theta_i) \quad (1)$$

and

$$\sigma^2 = \frac{\phi}{w_i} b''(\theta_i). \quad (2)$$

We assume that

$$\mu_i = h(\eta_i). \quad (3) \quad \left. \begin{array}{l} \text{possibly non-canonical} \\ \text{link function.} \end{array} \right\}$$

Combining (1) and (3)

$$b'(\theta_i) = h(\eta_i). \quad (4)$$

Implicit differentiation leads to

$$\frac{\partial}{\partial \beta} b'(\theta_i) = \frac{\partial}{\partial \beta} h(\eta_i)$$

$$b''(\theta_i) \frac{\partial \theta_i}{\partial \beta} = h'(\eta_i) \frac{\partial \eta_i}{\partial \beta}$$

such that

$$\frac{\partial \theta_i}{\partial \beta} = \frac{h'(\eta_i)}{b''(\theta_i)} \frac{\partial \eta_i}{\partial \beta} = \frac{h'(\eta_i)}{b''(\theta_i)} x_i$$

$$\eta_i = \tilde{x}_i \beta$$

$$\frac{\partial \eta_i}{\partial \beta} = x_i$$

Contribution to score vector

$$s_i(\beta) = \frac{w_i}{\phi} \left( y_i \frac{\partial \theta_i}{\partial \beta} - b'(\theta_i) \frac{\partial \theta_i}{\partial \beta} \right)$$

$$= \frac{w_i}{\phi} (y_i - \mu_i) \frac{h'(y_i)}{b''(\theta_i)} \underbrace{x_i}_{\substack{n \\ \text{~} \\ \text{~} \\ \text{~} \\ p \times 1}}$$

$$= x_i (y_i - \mu_i) \frac{h'(y_i)}{\sigma_i^2}$$

Score vector

$$s(\beta) = \sum_{i=1}^n s_i(\beta) = \underbrace{X^\top D}_{p \times n} \sum^{-1} (\underbrace{y - \mu}_{n \times 1})$$

where  $D = \text{diag}(h'(y)) = \begin{bmatrix} h'(y_1) & & & \\ & h'(y_2) & & \\ & & \ddots & \\ & & & h'(y_n) \end{bmatrix}$  and  $\sum = \text{diag}\left(\frac{1}{\sigma_i^2}\right) = \begin{bmatrix} \frac{1}{\sigma_1^2} & & & \\ & \frac{1}{\sigma_2^2} & & \\ & & \ddots & \\ & & & \frac{1}{\sigma_n^2} \end{bmatrix}$

Fisher information:

$$\begin{aligned} F(\beta) &= \underbrace{\sum_{i=1}^n \text{Var}(s_i(\beta))}_{p \times p} \\ &= \sum_{i=1}^n x_i^\top \text{Var}\left(\frac{h'(y_i)}{\sigma_i^2} y_i\right) x_i \\ &= \sum_{i=1}^n x_i^\top \frac{(h'(y_i))^2}{\sigma_i^2} \end{aligned}$$

$$= X^\top W X \quad \text{where} \quad W = D^2 \sum^{-1}$$

Fisher scoring algorithm

$$\begin{aligned}\hat{\beta}_{t+1} &= \hat{\beta}_t + F(\hat{\beta}_t)^{-1} \Sigma(\hat{\beta}_t) \\ &= (\hat{X}^T W \hat{X})^{-1} \hat{X}^T W X \hat{\beta}_t + (\hat{X}^T W \hat{X})^{-1} \hat{X}^T D \Sigma^{-1} (\hat{y} - \hat{\mu}) \\ &= (\hat{X}^T W \hat{X})^{-1} \hat{X}^T W (\hat{X} \hat{\beta}_t + D^{-1} (\hat{y} - \hat{\mu})) \\ &= (\hat{X}^T W \hat{X})^{-1} \hat{X}^T W \hat{y} \quad (\hat{\beta}_{MLE} \text{ for general linear model})\end{aligned}$$

where the "working observations"

$$\hat{y}_i^{(t)} = y_i^{(t)} + \frac{y_i - h(y_i^{(t)})}{h'(y_i^{(t)})}$$

## Iterated Reweighted Least Squares (IRLS)

Advantages:

- Instead of initial values  $\hat{\mu}$   $glm(\ , \text{start})$   
we can instead initialise  $\hat{\mu} = h(\hat{y})$   $glm(\ , \text{mustart=}\cdot)$

By default, initial values for  $\hat{\mu}$  are computed by the \$initialise function component of the family argument, returning

$$\hat{\mu}_i = \frac{y_i + 0.5}{n_i + 1}$$

for binomial GLMs.

- Builds on efficient code computing

$$\hat{\beta}_{t+1} = (\hat{X}^T W \hat{X})^{-1} W \hat{y}$$

via the QR-decomposition of  $W^{1/2} X$   
(e.g. via Gram-Schmidt process,

more stable numerically)

instead of via inversion of  $X^T W X$   
(numerically unstable)

Both have computational complexity  $O(n p^2)$

General linear model via QR-decomp:

$$\underline{y} = X\beta + \underline{\varepsilon}, \quad \underline{\varepsilon} \sim N(\underline{0}, \Sigma), \quad \Sigma^{-1} = L^T L$$

Then

$$\underline{\varepsilon} = (L^T L)^{-1}$$

$$Ly = L X \beta + L \underline{\varepsilon}$$

$$\text{or } \underline{y}^* = X^* \hat{\beta} + \hat{\underline{\varepsilon}}$$

where

$$\text{Var}(\hat{\underline{\varepsilon}}) = \text{Var}(L \underline{\varepsilon}) = L \Sigma L^T = L (L^T L)^{-1} L^T = L L^{-1} L^T L^T = I$$

Hence,  $\hat{\beta}$  satisfies

$$X^{*\top} (\underline{y}^* - X^* \hat{\beta}) = 0 \quad (\text{columns of } X^* \text{ orthogonal to } \hat{\underline{\varepsilon}})$$

Let

$$QR = X^*$$

denote QR-decomposition of  $X^*$  ( $Q^T Q = 0$  and R upper triangular) (Gram-Schmidt)

Then

$$(QR)^T (\underline{y} - QR \hat{\beta}) = 0$$

$$R^T Q^T \underline{y} = R^T \underbrace{Q^T Q}_{=I} R \hat{\beta}$$

$$R \hat{\beta} = Q^T \underline{y}$$

## Quasi-likelihood (ch. 5.5)

For real likelihoods from exponential family

$$s_i(\beta) = \underline{x}_i \frac{h'(y_i)}{\sigma^2(\mu_i)} (y_i - \mu_i) \quad (*)$$

and

$$F_i(\beta) = \underline{x}_i \underline{x}_i^\top \frac{(h'(y_i))^2}{\sigma^2(\mu_i)}$$

depends on  $f(y_i | \mu_i, \dots)$  only via  $\mu_i$  and  $\sigma_i^2 = \sigma^2(\mu_i)$

Idea: Accommodate overdispersion (and other relationships between  $\sigma_i^2$  and  $\mu_i^2$ ) by changing  $\sigma_i^2 \rightarrow \sigma_i^2 = \phi \sigma^2(\mu_i)$ . Does not necessarily correspond to any real underlying distribution.

Quasi-binomial:

$$s_i(\beta) = \underline{x}_i \frac{h'(y_i)}{\phi \pi_i(1-\pi_i)} n_i (\bar{y}_i - \pi_i)$$

$$F_i(\beta) = \underline{x}_i \underline{x}_i^\top \frac{n_i (h'(y_i))^2}{\phi \pi_i(1-\pi_i)}$$

Quasi-Poisson

$$s_i(\beta) = \underline{x}_i \frac{h'(y_i)}{\phi \lambda_i} (y_i - \lambda_i)$$

$$F_i(\beta) = \underline{x}_i \underline{x}_i^\top \frac{(h'(y_i))^2}{\phi \lambda_i}$$

Note: The MLE of  $\beta$  remain unchanged but  $\text{Var}(\hat{\beta})$  is inflated by a factor  $\phi$ .

Estimate of  $\phi$

$$\hat{\phi} = \frac{1}{n-p} \sum_{i=1}^n \frac{y_i - \hat{\mu}_i}{\sqrt{\sigma^2(\hat{\mu}_i)}} \quad \text{or alternatively} \quad \hat{\phi} = \frac{D}{n-p}$$

R:  $\text{glm}(y \sim x, \text{family} = \text{quasipoisson}(\text{link} = \text{"log"}))$

Corresponds to a quasi-likelihood with contributions defined as

$$Q_i(\mu_i; y_i) = \int_{y_i}^{\mu_i} \frac{y_i - t}{\phi \sigma^2(t)} dt \quad (\text{exercise 11})$$

Note how  $\frac{\partial Q}{\partial \mu}$  equals (\*). We define quasi-deviance contribution as

$$D_i(\hat{\mu}_i; y_i) = -2\phi [Q_i(y_i; y_i) - Q_i(\hat{\mu}_i; y_i)]$$

$$= +2\phi Q_i(\hat{\mu}_i; y_i)$$

quasibinomial

$$= y_i \ln \frac{\hat{\mu}_i}{\bar{\mu}_i} + (n_i - y_i) \ln \frac{1 - \hat{\mu}_i}{1 - \bar{\mu}_i}$$

Given two nested models  $H_0$  and  $H_1$

$$\frac{D_0 - D_1}{\phi} \stackrel{\text{asym}}{\sim} \chi^2_{p_1 - p_0}$$

Also

$$\frac{D_1}{\phi} \stackrel{\text{asym}}{\sim} \chi^2_{n - p_1}$$

Under  $H_0$

$$F = \frac{(D_1 - D_0) / (\rho_1 - \rho_0)}{\hat{\phi}}$$

$$= \frac{\left( \frac{D_1 - D_0}{\hat{\phi}} \right) / (\rho_1 - \rho_0)}{\frac{D_1}{\hat{\phi}} / (n - \rho_1)} \stackrel{\text{asympt}}{\sim} F_{\rho_1 - \rho_0, n - \rho_1}$$

see Venables & Ripley, eq. 7.10. (springer link)

R: `drop1(mod, test = "F")`

`anova(mod0, mod1, test = "F")`

Replaces `test = "LRT"` for non-quasi-likelihood models.

Alternative approaches:

Parametric models:

Beta-binomial:  $y_i | p_i \sim \text{bin}(n_i, p_i)$

$p_i \sim \text{Beta}(\alpha, \beta)$

Negative binomial:  $y_i | \lambda_i \sim \text{Poisson}(\lambda_i)$

$\lambda_i \sim \text{Gamma}(\alpha, \beta)$

GLMMs

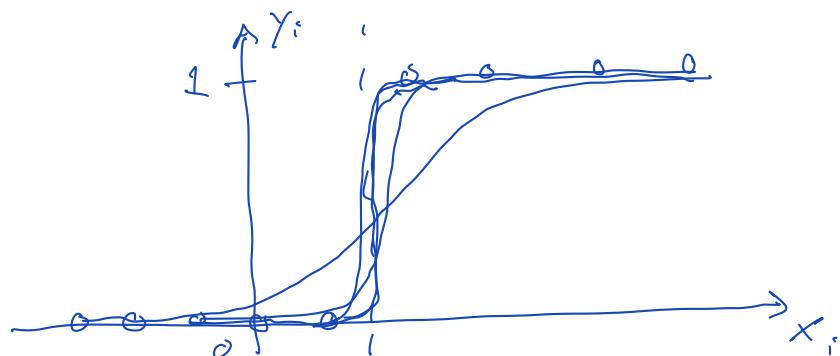
Models with random variables in  
the linear predictor

## Linear separation

Sometimes encountered when working with sparse data, typically if there are a small number of observations for some levels of a factor.

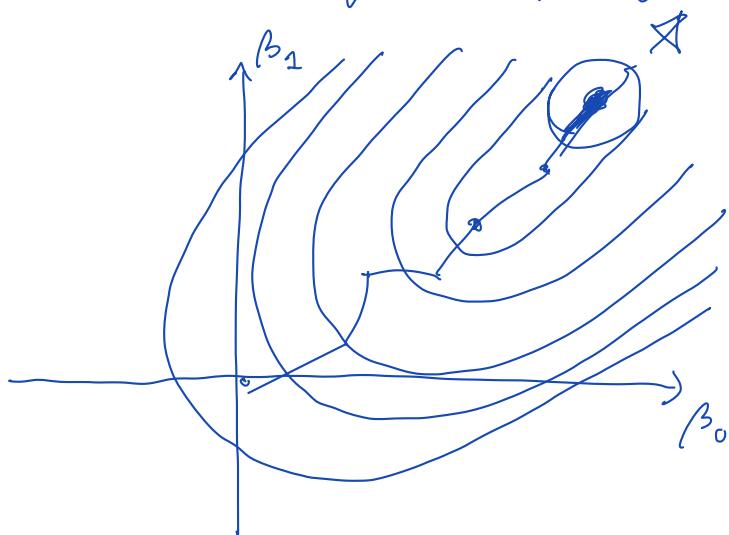
Minimal example:  $y_i \sim \text{bin}(1, \pi_i)$ ,

$$\log(\pi_i) = \beta_0 + \beta_1 x_i$$



The MLEs of  $\beta_0$  and  $\beta_1$  goes to infinity.  $F(\beta)$  will

have one eigenvalue very close to zero and  $\text{Var}(\hat{\beta}) = F^{-1}(\beta)$  will become extremely large.



If we do a Wald-test based on  $\text{Var}(\hat{\beta})$

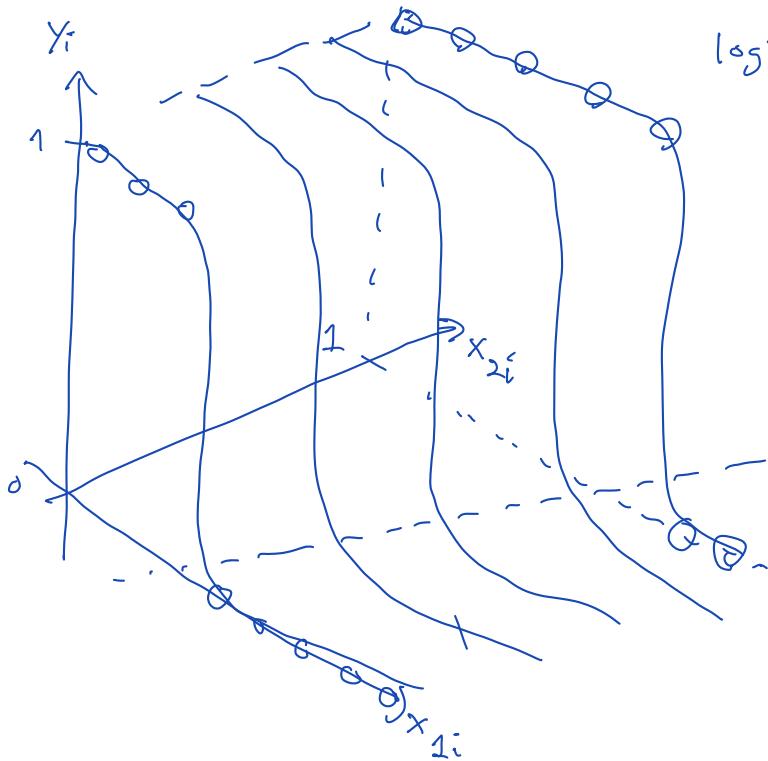
$$Z = \frac{\hat{\beta}_1}{\text{SE}(\hat{\beta}_1)}$$

becomes very small.  
The Wald test can't be trusted.

$H_0: \beta_1 = 0$  vs  $H_1: \beta_1 \neq 0$  can still be tested

using  $LRT = 2(\ell(\hat{\beta}_1) - \ell(\hat{\beta}_0)) \sim \chi^2_{p_1 - p_0}$

Example 2: Two covariates



$$\text{logit } \pi_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}$$

No linear separation  
in terms of  $x_{1i}$  or  
 $x_{2i}$  alone.

A certain linear combination  $L = a_1 x_{1i} + a_2 x_{2i}$   
separates  $y_i$  into 0's and 1's.

Possible solutions:

Penalized regression: Lasso, ridge

Using GLMMs modelling effect associated  
with each level as random  $\underline{x}_i \sim N(0, \Sigma)$

## Offset variables

Additional term in linear predictor involving no unknown parameters

$$y_i = \hat{x}_i^\top \beta + \alpha_i$$

R:  $\text{glm}(y \sim x, \text{offset} = \alpha_1, \dots)$

1) Useful for likelihood ratio test of linear hypotheses

$$C\beta = d, \text{ e.g. } H_0: \beta_1 = 1 \text{ vs } H_1: \beta_1 \neq 1$$

e.g. the Poisson regression model

$$y_i \sim \text{Poisson}(\lambda_i), \lambda_i = c x_i^\beta$$

$$\ln \lambda_i = \beta_0 + \beta_1 \ln x_i$$

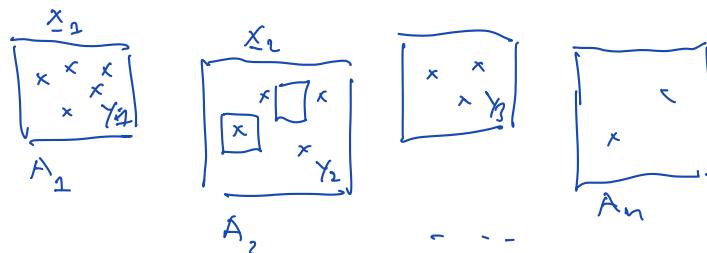
$H_0: \text{glm}(y \sim 1, \text{offset} = \log(x), \text{poisson})$

$H_1: \text{glm}(y \sim \log(x), \text{poisson})$

`anova(mod0, mod1, test = "Chisq")`

(Also see exercise 2)

2) Poisson regression and rates



We expect  $E(y_i) = A_i e^{\hat{x}_i^\top \beta}$

Model:  $\lambda_i = A_i e^{\hat{x}_i^\top \beta}$

$$\ln \lambda_i = \hat{x}_i^T \beta + \underbrace{\ln A_i}_{\text{offset term}}$$

Confidence intervals for regression coefficients or functions of  $\beta$  more generally

- Wald: We have

$$P\left(-z_{\alpha/2} < \frac{\hat{\beta}_j - \beta_j}{\widehat{SE}(\hat{\beta}_j)} < z_{\alpha/2}\right) \approx 1 - \alpha$$

↓  
 asymp  
 $\sim N(0, 1)$   
 pivotal quantity

Hence,

$$P\left(\hat{\beta}_j - \widehat{SE}(\hat{\beta}_j) z_{\alpha/2} < \beta_j < \hat{\beta}_j + \widehat{SE}(\hat{\beta}_j) z_{\alpha/2}\right) \approx 1 - \alpha$$

Thus

$$\hat{\beta}_j \pm \widehat{SE}(\hat{\beta}_j) z_{\alpha/2}$$

is a  $(1-\alpha)$ -conf. int. for  $\beta_j$

- Profile likelihood confidence intervals

Recall: For a given value  $\beta_j$  of the  $j$ th component of  $\beta$

$$P\left(2\left(l(\hat{\beta}_j, \hat{\beta}_{-j}) - l(\beta_j, \hat{\beta}_{-j,0})\right) < \chi^2_{\alpha, 1}\right) \approx 1 - \alpha$$

↓  
 pivotal quantity

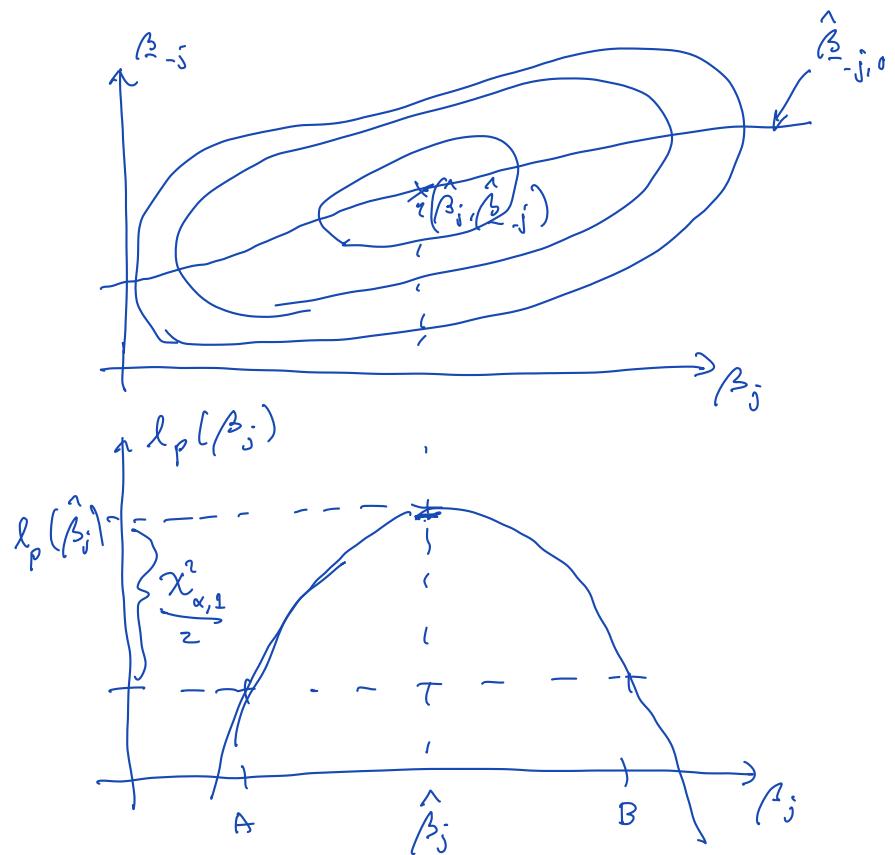
Letting  $l_p(\beta_j) = \sup_{\beta_{-j}} l(\beta_j, \beta_{-j})$  denote the profile likelihood for  $\beta_j$  ("profiling out"  $\beta_{-j}$ )

(\*) can be written

$$P\left(2\left(l_p(\hat{\beta}_j) - l_p(\beta_j)\right) < \chi^2_{\alpha, 1}\right) \approx 1 - \alpha$$

$$P\left(l_p(\beta_j) > l_p(\hat{\beta}_j) - \frac{\chi^2_{\alpha, 1}}{2}\right) \approx 1 - \alpha$$

$$P(A < \beta_j < B) \approx 1-\alpha$$



R::confint(mod) calls MASS:::confint.glm  
 $\text{mle}(\text{fit})$  with  $\beta_j x_{ij}$  included  
as an offset term to compute  $\ell_p(\beta_j)$

Functional invariance:

If  $(A, B)$  is the profile conf.int for  $\theta$ , then  $(g(A), g(B))$  is the profile likelihood conf.int. for  $g(\theta)$   
provided that  $g$  is strictly monotonic.

## Vector GLM (multinomial regression Ch. 6)

More than two categories  $\{1, 2, \dots, c, c+1\}$

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_c \end{bmatrix} \sim \text{Multinomial}(m, \underline{\pi})$$

$m$   
 $c \times 1$

Pmf of multinomial distn.

$$f(\underline{y} | \underline{\pi}) = \frac{m!}{y_1! y_2! \dots y_c! (m-y_1-\dots-y_c)!} \pi_1^{y_1} \pi_2^{y_2} \dots (1-\pi_1-\dots-\pi_c)^{m-y_1-\dots-y_c}$$

$$\mathbb{E}(\underline{y}) = m \underline{\pi}$$

$$\text{Var}(\underline{y}) = m \underbrace{\begin{bmatrix} \pi_1(1-\pi_2) & & & \\ -\pi_1\pi_2 & \pi_2(1-\pi_3) & & \\ & -\pi_2\pi_3 & \ddots & \\ & & -\pi_{c-1}\pi_c & \pi_c(1-\pi_c) \end{bmatrix}}_{c \times c}$$

Data:

Response matrix

$$\begin{bmatrix} \underline{y}_1^T \\ \underline{y}_2^T \\ \vdots \\ \underline{y}_n^T \end{bmatrix} = \begin{bmatrix} y_{11} & y_{12} & \dots & y_{1c} \\ y_{21} & & & \\ \vdots & & & \\ y_{n1} & y_{n2} & \dots & \end{bmatrix}$$

$n \times c$

Design matrix

$$\begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots \\ x_{21} & & \\ \vdots & & \\ \vdots & & \end{bmatrix}$$

$n \times p$

## 6.2. Models for unordered data (nominal multinomial regression)

Tentative model: Let  $\pi_{ir} \propto e^{\underline{x}_i^T \underline{\beta}_r}$  and normalize, i.e.

$$\pi_{ir} = \frac{e^{\underline{x}_i^T \underline{\beta}_r}}{\sum_{s=1}^{c+1} e^{\underline{x}_i^T \underline{\beta}_s}}$$

} specific to each category r.

$p \times (c+1)$  parameters

Adding a constant vector  $\underline{f}$  to all the  $\underline{\beta}_r$ 's doesn't change the  $\pi_{ir}$  (and the distr. of all the data) since

$$\pi_{ir}^* = \frac{e^{\underline{x}_i^T (\underline{\beta}_r + \underline{f})}}{\sum_{s=1}^{c+1} e^{\underline{x}_i^T (\underline{\beta}_s + \underline{f})}} = \pi_{ir}$$

Hence, the model is non-identifiable.

Identifiable model: Work with  $\underline{\beta}_r = \underline{\beta}_r - \underline{\beta}_{c+1}$ . This leads to

$$\pi_{ir} = \frac{e^{\underline{x}_i^T \underline{\beta}_r}}{1 + \sum_{s=1}^c e^{\underline{x}_i^T \underline{\beta}_s}} \quad \text{for } r = 1, 2, \dots, c$$

$$\pi_{i,c+1} = \frac{1}{1 + \sum_{s=1}^c e^{\underline{x}_i^T \underline{\beta}_s}} \quad \text{for } r = c+1$$

}  $p \times c$  parameters

- Odds ratio interpretation: A Unit change in  $x_{ij}$ , changing  $\underline{x}_i$  to  $\underline{x}_i^* = \underline{x}_i + (0, 0, 1, 0, \dots, 0)^T$  changes the odds element  $j$

of event  $r_1$  relative to  $r_2$  by a

factor (an odds ratio) of

$$\frac{\pi_{ir_1}^* / \pi_{ir_2}^*}{\pi_{ir_1} / \pi_{ir_2}} = \frac{e^{x_i^T (\beta_{r_1} - \beta_{r_2})}}{e^{x_i^T (\beta_{r_1} - \beta_{r_2})}} = e^{\beta_{r_1,j} - \beta_{r_2,j}}$$

$$= e$$

Example: Diet of alligators: 219 alligators  
 Response is diet (fish, invertebrates, reptiles,  
 bird, other)

Two covariates

Size: numerical (see r-code)

Lake: 4 levels

Likelihood

$$L(\beta) = \prod_{i=1}^n f(y_i | n_i, \pi_i)$$

depends on  $\beta$

$\downarrow$   
 $\underbrace{\phantom{f(y_i | n_i, \pi_i)}$   
 $c \times 1$

$c_p \times 1$

Log-likelihood

$$l(\beta) = \sum_{i=1}^n \sum_{s=1}^{c+1} y_{is} \ln(\pi_{is}(\beta))$$

Deviance

$$D = 2 \sum_{i=1}^n \sum_{s=1}^{c+1} y_{is} \ln \left( \frac{y_{is} / n_i}{\hat{\pi}_{is}} \right)$$

asympt/ approx  $\sim \chi^2_{nc - pc}$

$\underbrace{nc}_{p_1} - \underbrace{pc}_{p_0}$

LRTs / test based on change in deviance work as before.

Goodness-of-fit via observed deviance

→ Latent utility interpretation of the model

Assume that individuals ("consumers") assess / perceive the utility of different alternatives ("products")  $r = 1, 2, \dots, c+1$  with some error

$$u_r = \gamma_r + \varepsilon_r$$

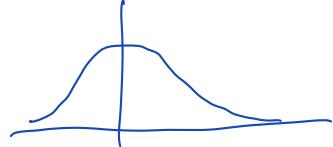
$\sim$

$\stackrel{\text{iid}}{\sim} \text{Gumbel}(0, 1)$

$$F_{\varepsilon_r}(x) = e^{-e^{-x}}$$

$$\text{Gumbel}(\mu, \sigma)$$

$$F(x) = e^{-e^{-\frac{x-\mu}{\sigma}}}$$



and choose the alternative with the largest  $u_r$ .

Theorem (McFadden 1973): It then follows that

$$P(\text{choosing } r) = \frac{e^{\gamma_r}}{\sum_{s=1}^{c+1} e^{\gamma_s}}.$$

Proof:

$$P(\text{choosing } r) = P\left(u_r > \max_{s \neq r} u_s\right)$$

$$= P \left( u_r - \max_{s \neq r} u_s > 0 \right)$$

$\sim$   $\text{Gumbel}(\gamma_r, 1)$        $\sim ?$

Cdf of the maximum:

$$\begin{aligned}
 F_{\max_{s \neq r} u_s}(x) &= P \left( \max_{s \neq r} u_s < x \right) \\
 &= P \left( \bigcap_{s \neq r} u_s < x \right) \\
 &= \prod_{s \neq r} P(u_s < x) \\
 &= \prod_{s \neq r} e^{-(x - \eta_s)} \\
 &= \prod_{s \neq r} e^{-e^{-x} \sum_{s \neq r} e^{\eta_s}} \\
 &= e^{-e^{-x} \sum_{s \neq r} e^{\eta_s}} \\
 &= e^{-e^{-x} \underbrace{\ln \sum_{s \neq r} e^{\eta_s}}_{\mu}}
 \end{aligned}$$

That is, Gumbel  $\left( \ln \sum_{s \neq r} e^{\eta_s}, 1 \right)$

If  $X \sim \text{Gumbel}(\mu, \sigma)$  then

$$F(x) = e^{-e^{-x}}$$

$$= e^{-x - e^{-x}}$$

$$f(x) = \dots = e^{-x - e^{-x}}$$

$$M_X(t) = \int_{-\infty}^{\infty} e^{tx - x - e^{-x}} dx$$

$$= \int_{-\infty}^0 e^{-(t-1)\ln u - u} \left(-\frac{1}{u}\right) du$$

$$= \int_0^{\infty} u^{-t+1} e^{-u} du = \int_0^{\infty} u^{1-t-1} e^{-u} du$$

$$= \Gamma(1-t)$$

Gumbel( $\mu, \sigma$ ):

$$M_{\mu+\sigma X}(t) = E e^{t(\mu+\sigma X)} = e^{\mu t} \underbrace{E e^{t\sigma X}}$$

$$= e^{\mu t} M_X(t\sigma)$$

$$= e^{\mu t} \Gamma(1 - \sigma t)$$

Hence,

$$M_{u_r - \max_{s \neq r} u_s}(t) = e^{\gamma_r t} \Gamma(1 - t) e^{\left(\ln \sum_{s \neq r} e^{\gamma_s}\right)t} \Gamma(1 + t)$$

$$= e^{\left(\gamma_r - \ln \sum_{s \neq r} e^{\gamma_s}\right)t} \underbrace{\Gamma(1 - t) \Gamma(1 + t)}_{= B(1+t, 1-t)}$$

If  $X \sim \text{logistic}(0, 1)$  then

$$F(x) = \frac{e^x}{1+e^x}, \quad f(x) = \frac{e^x}{(1+e^x)^2}$$

$$F(x) = \frac{e^{x-\mu}}{1+e^{x-\mu}}$$

$$M_X(t) = \int_{-\infty}^{\infty} e^{tx} \frac{e^x}{(1+e^x)^2} dx$$

$$u = \frac{e^x}{1+e^x} \quad x = \ln \frac{u}{1-u}$$

$$= \int_0^1 u(1-u) \left(\frac{u}{1-u}\right)^t \frac{du}{u(1-u)}$$

$$dx = \frac{1-u}{u} \cdot \frac{1-u+u}{(1-u)^2} du$$

$$= \frac{du}{u(1-u)}$$

$$= B(1+t, 1-t)$$

Logistic( $\mu, \sigma$ ):

$$\begin{aligned} M_{\mu + \sigma X}(t) &= E e^{t(\mu + \sigma X)} \\ &= e^{t\mu} E(e^{t\sigma X}) \\ &= e^{t\mu} M_X(t\sigma) \end{aligned}$$

$$= e^{t\mu} B(1+\sigma t, 1-\sigma t)$$

Thus,  $u_r - \max_{s \neq r} u_s \sim \text{logistic}\left(y_r - \ln \sum_{s \neq r} e^{y_s}, 1\right)$

and

$$P(\text{Choosing } r) = P\left(u_r - \max_{s \neq r} u_s > 0\right)$$

$$= 1 - F_{\text{logistic}(\cdot, 1)}(0)$$

$$= 1 - \frac{e^{0 - \gamma_r + \ln \sum_{s \neq r} \gamma_s}}{1 + e^{0 - \gamma_r + \ln \sum_{s \neq r} \gamma_s}}$$

$$= 1 - \frac{\sum_{s \neq r} e^{\gamma_s} e^{-\gamma_r}}{1 + \sum_{s \neq r} e^{\gamma_s} e^{-\gamma_r}}$$

$$= 1 - \frac{\sum_{s \neq r} e^{\gamma_s}}{e^{\gamma_r} + \sum_{s \neq r} e^{\gamma_s}}$$

$$= \frac{e^{\gamma_r}}{\sum_{s=1}^{c+1} e^{\gamma_s}} \quad (\text{proving the theorem})$$

$$= \frac{e^{\sum_{s=1}^{c+1} \hat{\beta}_s x_i^s}}{\sum_{s=1}^{c+1} e^{\sum_{s=1}^{c+1} \hat{\beta}_s x_i^s}} \quad (\text{as in multinomial logit model})$$

## Ordinal regression (6.3)

The response is an ordered categorical variable  $Y_i \in \{1, 2, \dots, c+1\}$  e.g. grades (A, B, ..., F), a Likert scale (strongly disagree, disagree, indifferent, ..., strongly agree).

Model: The response  $Y_i$  is determined by an underlying latent variable

$$u_i = -x_i^\top \beta + \varepsilon_i$$

where  $\varepsilon_i$  has a standard distribution with cdf  $F$ .

The event  $Y_i = r$  occurs if

$$\theta_{r-1} < u_i \leq \theta_r$$

$$\text{where } -\infty = \theta_0 < \theta_1 < \dots < \theta_{c+1} = \infty$$

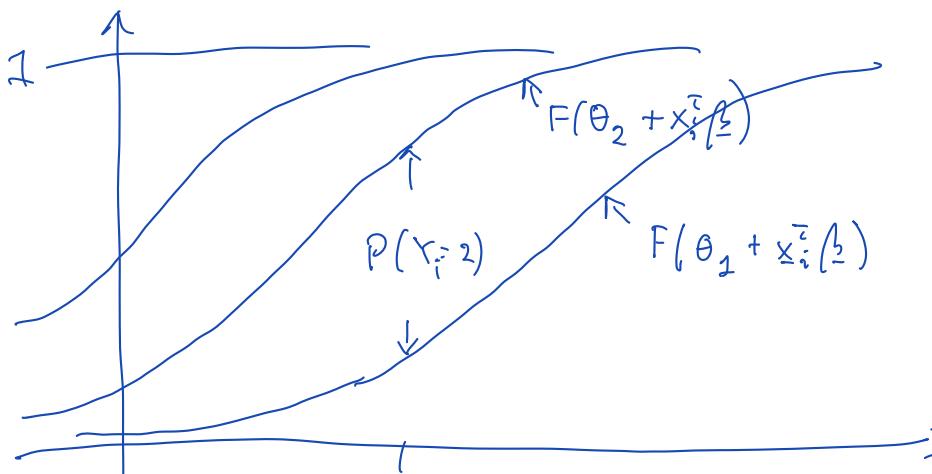
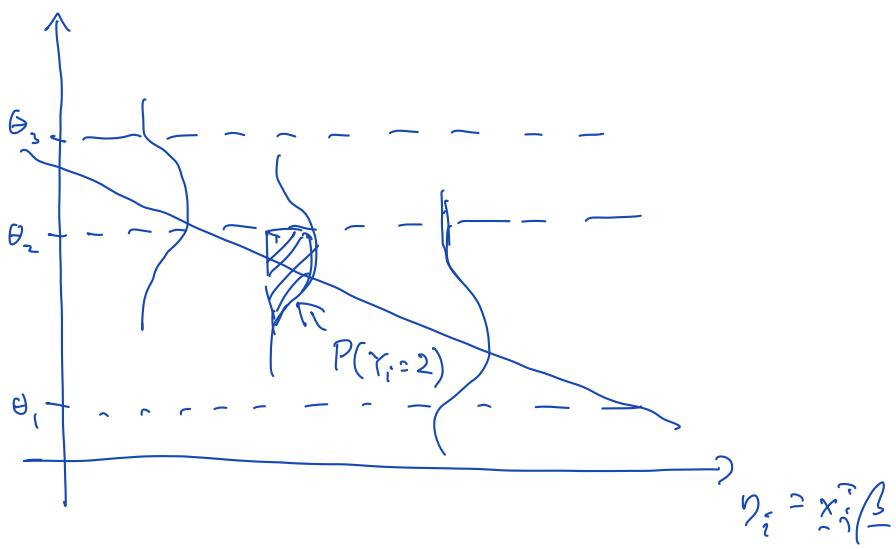
Unknown parameters,  $\theta_1, \theta_2, \dots, \theta_c, \beta_1, \beta_2, \dots, \beta_k$ .

Then

$$\begin{aligned} P(Y_i \leq r) &= P(u_i \leq \theta_r) \\ &= P(-x_i^\top \beta + \varepsilon_i \leq \theta_r) \\ &= P(\varepsilon_i \leq \theta_r + x_i^\top \beta) \\ &= F(\theta_r + x_i^\top \beta) \end{aligned}$$

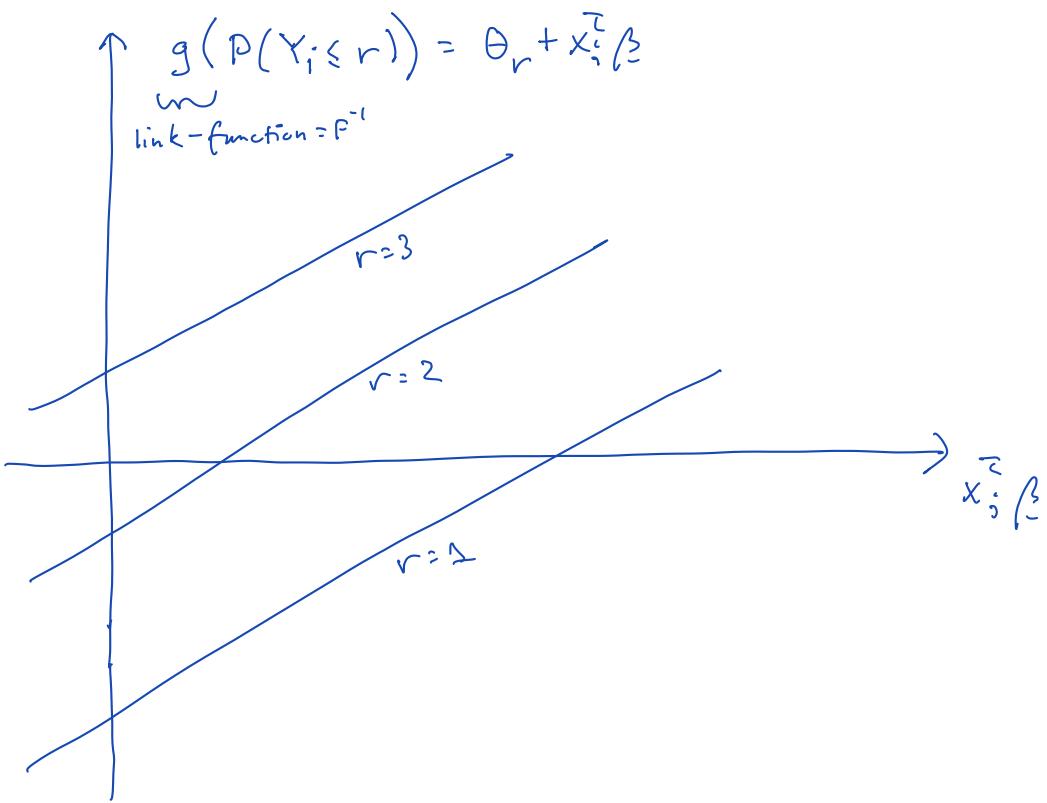
and

$$\begin{aligned} \pi_{ir} &= P(Y_i = r) = P(Y_i \leq r) - P(Y_i \leq r-1) \\ &= F(\theta_r + x_i^\top \beta) - F(\theta_{r-1} + x_i^\top \beta) \end{aligned}$$



$\beta > 0 \Rightarrow E(Y_i)$  decreases with  $x_i$

Parallelism assumption:



Proportional odds model:

For  $F(x) = \frac{e^x}{1+e^x}$ , that is,  $\varepsilon_i \sim \text{logistic}(0, 1)$

$$\theta_r + x_i^\top \beta$$

$$P(Y_i \leq r) = F(\theta_r + x_i^\top \beta) = \frac{e^{\theta_r + x_i^\top \beta}}{1 + e^{\theta_r + x_i^\top \beta}}$$

$$\text{Odds}(Y_i \leq r) = \frac{P(Y_i \leq r)}{P(Y_i > r)} = \dots = e^{\theta_r + x_i^\top \beta} = e^{\theta_r} e^{x_i^\top \beta}$$

Changing covariate  $x_{ij}$  by one unit changes all the odds of the cumulative events  $Y_i \leq 1, Y_i \leq 2, \dots$

by a common factor  $e^{\beta_j}$ , i.e. the odds of those events change proportionally.

## Cumulative probit

$F(x) = \Phi(x)$  perhaps a more natural assumption  
(discretisation of latent LM)

R: library(VGAM)

`vglm(y ~ . . ., family = cumulative(link = "probit",  
parallel = TRUE))`

↑  
`factor(, ordered = TRUE)`

## Likelihood inference (nominal and ordinal multinomial models) 6.4

General notation:

$$\underline{y}_i \sim \text{Multinomial}(n_i, \underline{\pi}_i)$$

$\underbrace{\phantom{y_i}}_{m}$        $\underbrace{\phantom{\pi_i}}_{m}$   
 $m$        $c \times 1$   
 $c \times 1$

Vector response function linking  $\underline{\pi}_i$  to  $\underline{y}_i$ ,

$$\underline{\pi}_i = h(\underline{\eta}_i)$$

$\underbrace{\phantom{\pi_i}}_m$        $\underbrace{\phantom{\eta_i}}_m$   
 $m$        $c \times 1$   
 $c \times 1$

where

$$\underline{\eta}_i = X_i \underline{\beta}.$$

$\underbrace{\phantom{\eta_i}}_{c \times 1}$        $\underbrace{\phantom{\beta}}_{c \times p}$   
 $c \times 1$        $c \times p$

Example: Nominal model with  $k$  covariates has response function

$$\underline{\pi}_i = h(\underline{\eta}_i) = \frac{\begin{bmatrix} e^{\eta_{i1}} \\ e^{\eta_{i2}} \\ \vdots \\ e^{\eta_{ic}} \end{bmatrix}}{\left(1 + \sum_{s=1}^c e^{\eta_{is}}\right)}$$

and

$$\underline{\eta}_i = \begin{bmatrix} \eta_{i1} \\ \eta_{i2} \\ \vdots \\ \eta_{ic} \end{bmatrix} = \begin{bmatrix} \underline{x}_i^\top & 0 & \cdots & 0 \\ 0 & \underline{x}_i^\top & & \\ \vdots & & \ddots & \\ 0 & & & \underline{x}_i^\top \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_c \end{bmatrix}$$

$c \times 1$

$c \times (k+1)c$

$(k+1)c \times 1$

Example: Ordinal model:

$$\underline{\pi}_i = h(\underline{\eta}_i) = \begin{bmatrix} F(\eta_{i1}) - 0 \\ F(\eta_{i2}) - F(\eta_{i1}) \\ \vdots \\ 1 - F(\eta_{ic}) \end{bmatrix}$$

$$\underline{\eta}_i = \begin{bmatrix} \eta_{i1} \\ \eta_{i2} \\ \vdots \\ \eta_{ic} \end{bmatrix} = \begin{bmatrix} 1 & 0 & \cdots & 0 & \underline{x}_i^\top \\ 0 & 1 & & & \vdots & \underline{x}_i^\top \\ \vdots & & 1 & \cdots & \vdots & \vdots \\ 0 & & \cdots & 1 & \vdots & \vdots \\ 0 & & & & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_c \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$$

$$c \times (c+k) \quad (c+k) \times 1$$

The log-likelihood contribution from obs.  $i$  becomes

$$l_i(\beta) = \sum_{s=1}^c y_{is} \ln \pi_{is} + (n_i - y_{i1} - \dots - y_{ic}) \ln (1 - \pi_{i2} - \dots - \pi_{ic})$$

This leads to the score contribution

$$\underline{s}_i(\beta) = \frac{\partial}{\partial \beta} l_i(\beta) = \dots = \underbrace{X_i^\top}_{p \times c} \underbrace{D_i}_{c \times c} \underbrace{\sum_i^{-1}}_{c \times c} \left( \underbrace{y_i - n_i \pi_i}_{c \times 1} \right)$$

where

$$\sum_i = \text{Var}(\underline{y}_i) = n_i \left( \text{diag}(\pi_i) - \underbrace{\pi_i \pi_i^\top}_{c \times c} \right)$$

Analogous to corresponding expression  
for GLMs

$$\underline{s}_i(\beta) \approx \underline{x}_i \frac{h'(x_i)}{\sigma_i^2} (y_i - \mu_i)$$

and

$$D_i = \frac{\partial h^\top}{\partial \underline{y}_i}$$

$$\underbrace{\phantom{D_i =}}_{c \times 1} \underbrace{\phantom{h^\top}}_{1 \times c}$$

Contribution to Fisher information

$$F_i(\beta) = \text{Var}(\underline{s}_i(\beta)) = \dots = \underbrace{X_i^\top}_{W_i} \underbrace{D_i \sum_i^{-1} D_i^\top}_{W_i^\top} X_i$$

Letting

$$\underline{Y} = \begin{bmatrix} \underline{Y}_1 \\ \underline{Y}_2 \\ \vdots \\ \underline{Y}_n \end{bmatrix}, \quad \underline{\Pi} = \begin{bmatrix} \underline{\Pi}_1 \\ \underline{\Pi}_2 \\ \vdots \\ \underline{\Pi}_n \end{bmatrix}, \quad N = \text{blockdiag}(n_1 \underline{\Gamma}_c, n_2 \underline{\Gamma}_c, \dots)$$

$$\underline{X} = \begin{bmatrix} \underline{X}_1 \\ \underline{X}_2 \\ \vdots \\ \underline{X}_n \end{bmatrix}, \quad \Sigma = \text{blockdiag}(\Sigma_1, \Sigma_2, \dots, \Sigma_n)$$
$$D = \text{blockdiag}(D_1, \dots)$$
$$W = \text{blockdiag}(W_1, \dots)$$

we can write

$$\underline{\zeta}(\underline{\beta}) = \sum_{i=1}^n \underline{\zeta}_i(\underline{\beta}) = \underline{X}^T D \Sigma^{-1} (\underline{Y} - N \underline{\Pi}),$$

and

$$F(\underline{\beta}) = \sum_{i=1}^n F_i(\underline{\beta}) = \underline{X}^T W \underline{X}.$$

Fisher Scoring algorithm:

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

# Ordinal regression of chess game outcomes (project 2)

Example data:

Tournament with  $k=3$  players,  $y_i = \begin{cases} 1 & \text{white win} \\ 2 & \text{draw} \\ 3 & \text{black win} \end{cases}$ ,  $n = \frac{k(k+1)}{2}$  games

Tentative model.

$$u_i = -x_i\beta + \varepsilon_i$$

$$-(\alpha_{j(i)} + \beta_{k(i)}) + \varepsilon_i \quad (*)$$

↑ effect  
of player  
 $j(i)$  having  
white pieces

↑ effect  
of player  
 $k(i)$  playing  
black

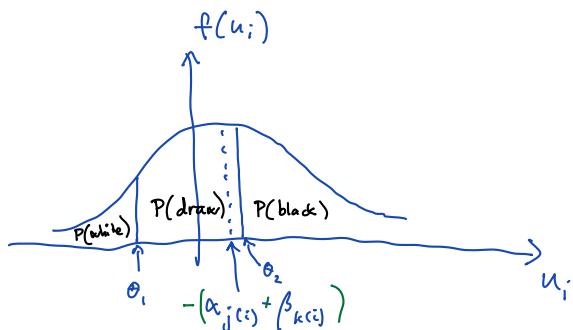
$\sim N(0, 1)$

Simpler model:

$$x_j = -\beta_j \quad \text{for } j=1, 2, \dots, k$$

$$u_i = \sim (\alpha_j - \beta_k) + \varepsilon_i \quad (**)$$

`vglm(y ~ x2 + x3, ...)`



`vglm(y ~ white + black, cumulative(...))`

| Game i | Outcome $y_i$ | player with white pieces $j(i)$ | player with black pieces $k(i)$ |
|--------|---------------|---------------------------------|---------------------------------|
| 1      | 1             | 1                               | 2                               |
| 2      | 1             | 1                               | 3                               |
| 3      | 2             | 2                               | 1                               |
| 4      | 2             | 2                               | 3                               |
| 5      | 1             | 3                               | 1                               |
| 6      | 1             | 3                               | 2                               |

Model matrix  
for (\*):

$$\left[ \begin{array}{cccccc} 1 & 0 & 0 & 1 & 0 & \vdots \\ 0 & 1 & 0 & 0 & 1 & \vdots \\ \vdots & \vdots & 1 & 0 & 0 & \alpha_2 \\ 0 & 1 & 1 & 0 & 0 & \alpha_3 \\ 0 & 0 & 1 & 0 & 0 & \beta_2 \\ 0 & 0 & 0 & 1 & 0 & \beta_3 \end{array} \right] \left[ \begin{array}{c} \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \beta_2 \\ \beta_3 \end{array} \right]$$

II

Model matrix  
for (\*\*):

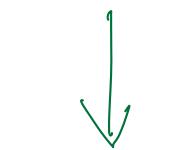
$$\left[ \begin{array}{ccccc} -1 & 0 & 0 & 0 & \vdots \\ 0 & -1 & 0 & 0 & \vdots \\ 1 & 0 & 0 & 0 & \alpha_1 \\ \vdots & \vdots & \vdots & \vdots & \alpha_2 \\ 1 & -1 & 0 & 0 & \alpha_3 \\ 0 & 1 & 0 & 0 & \beta_1 \\ -1 & 1 & 0 & 0 & \beta_2 \\ 0 & 0 & 1 & 0 & \beta_3 \end{array} \right] \left[ \begin{array}{c} \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{array} \right]$$

$x_2 \quad x_3$

`model.matrix(~ white + black)[-1]`

Does this matrix have  
full rank?

F. 60



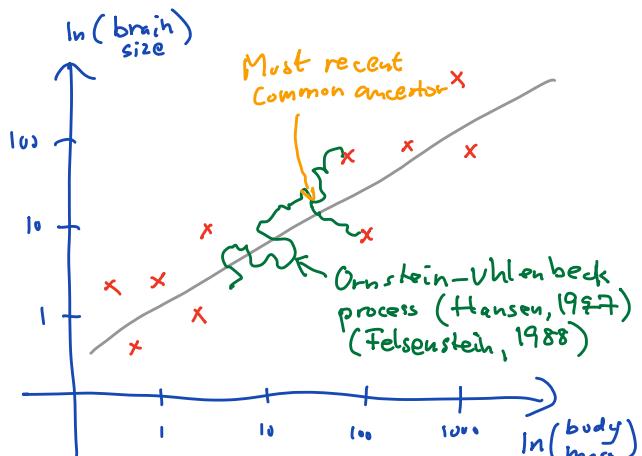
## Mixed models: LMMs and GLMMs (Ch 7)

↑  
19.10

General aim: Model and account for non-independent data, estimate dependency structure.

Examples:

- Repeated measurements on the same individuals / from the same group ("clustered data").
- Multiple grouping variables (crossed random effects)
- Temporal/spatial autocorrelation (time series analysis)
- Related individuals  (quantitative genetics)
- Related species



- Smooth: Effect  $z_t$  of age  $t$  given by a second order random walk

$$(1-\beta)^2 z_t = w_t$$

Notation: For clusters  $i = 1, 2, \dots, m$  we measure the response  $y_{ij}$  for different covariate vectors  $x_{ij}$  for  $j = 1, 2, \dots, n_i$

Random intercept model:

For  $i = 1, 2, \dots, m$ , and  $j = 1, 2, \dots, n_i$ :

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + \underbrace{\gamma_{0,i}}_{\text{random effect}} + \varepsilon_{ij}$$

where

$$\varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$$

and

$$\gamma_{0,i} \stackrel{iid}{\sim} N(0, \tau_0^2) \quad (*)$$

Alternatively: Model  $\gamma_{0,i}$  as a fixed effect?

Advantages of mixed model:

- We "borrow strength" between clusters.  
higher statistical power (project 3) at the cost of making one more assumption (\*)

R: library(lme4)

$$\text{lmer}(y \sim 1 + x + (1 | \underbrace{\text{group}}_{\text{grouping factor}}))$$

Conditional model: Within cluster  $i$  observations are indep.

$$Y_{ij} | \gamma_{0,i} \sim N(\beta_0 + \beta_1 x_{ij} + \gamma_{0,i}, \sigma^2)$$

Marginal model: Within clusters,

$$\begin{aligned} \text{Var}(y_{::j}) &= \text{Var}(\beta_0 + \beta_1 x_{::j} + \gamma_{0,:} + \varepsilon_{::j}) \\ &= \tau_0^2 + \sigma^2 \end{aligned}$$

$$\begin{aligned} \text{Cov}(y_{ij}, Y_{il}) &= \text{Cov}(\beta_0 + \beta_1 x_{ij} + \gamma_{0,i} + \varepsilon_{ij}, \\ &\quad \beta_0 + \beta_1 x_{il} + \gamma_{0,l} + \varepsilon_{il}) \\ &= \text{Cov}(\gamma_{0,i}, \gamma_{0,l}) = \text{Var}(\gamma_{0,:}) = \tau_0^2 \end{aligned}$$

## Intraclass correlation

$$\text{Corr}(Y_{ij}, Y_{ik}) = \frac{\tau_0^2}{\sigma^2 + \tau_0^2}$$

$$\underline{Y}_i = \begin{bmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{in_i} \end{bmatrix}_{n_i \times 2} \sim N\left(\underbrace{\underline{x}_i \beta}_{\text{mean}}, \sigma^2 I_{n_i} + \tau_0^2 J_{n_i}\right) \quad \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \end{bmatrix}$$

If incorrectly using a LM (no random effect)  
we still have

$$E(\hat{\beta}_{LM}) = E((\hat{X}^\top \hat{X})^{-1} \hat{X}^\top \hat{Y}) = \dots = \beta, \quad (\text{why?})$$

but

$$\text{Var}(\hat{\beta}_{LM}) > \sigma^2 (\hat{X}^\top \hat{X})^{-1}$$

with non-independent data. Also  $\hat{\beta}_{LM}$  not efficient.

## Random intercept and slope models (7.1.2)

For  $i=1, \dots, m$  and  $j=1, \dots, n_i$

$$Y_{ij} = \beta_0 + \beta_1 X_{ij} + \gamma_{0,i} + \gamma_{1,i} X_{ij} + \varepsilon_{ij}$$

$$= (\beta_0 + \gamma_{0,i}) + (\beta_1 + \gamma_{1,i}) X_{ij} + \varepsilon_{ij}$$

pop. mean  
intercept and slope

random  
deviations

where

$$\varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$$

and

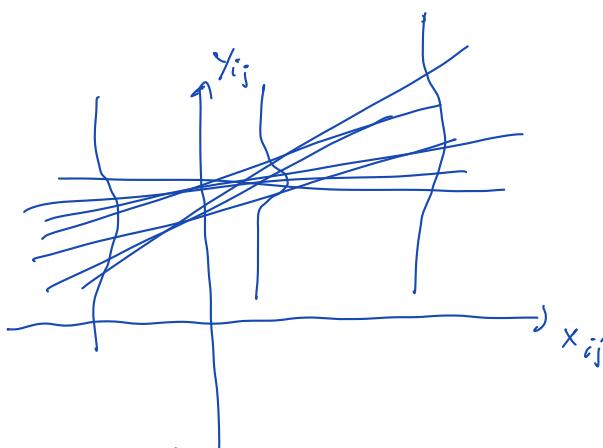
$$\underline{Y}_i = \begin{bmatrix} Y_{0,i} \\ Y_{1,i} \end{bmatrix} \stackrel{iid}{\sim} N \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \tau_0^2 & \tau_{01} \\ \tau_{01} & \tau_1^2 \end{bmatrix} \right)$$

$\underbrace{\qquad\qquad\qquad}_{Q}$

Marginal

$$\text{Var}(y_{ij}) = \tau_0^2 + \tau_1^2 x_{ij}^2 + 2x_{ij}\tau_{01} + \sigma^2$$

(quadratic in  $x_{ij}$ )



$$\begin{aligned} \text{Cov}(y_{ij}, y_{il}) &= \text{Cov} \left( \underbrace{Y_{0,i} + Y_{1,i}x_{ij} + \epsilon_{ij}}_{\text{fixed effect part}}, \underbrace{Y_{0,l} + Y_{1,l}x_{il} + \epsilon_{il}}_{\text{random intercept and slope}} \right) \\ &= \tau_0^2 + \tau_1^2 x_{ij} x_{il} + \tau_{01}(x_{ij} + x_{il}) \end{aligned}$$

R: lmer ( $y \sim 1 + x + (\underbrace{1 + x}_{\text{fixed effect part}} | \text{group})$ )  
 $\qquad\qquad\qquad \underbrace{\qquad\qquad\qquad}_{\text{random intercept and slope}}$

More general notation (7.1.3-7.1.4, 7.2)

Data  $(y_{ij}, \underline{x}_{ij}^\top)$  for  $i=1, 2, \dots, m$ ,  $j=1, 2, \dots, n_i$

$$y_{ij} = \underline{x}_{ij}^\top \underline{\beta} + u_{ij}^\top \underline{\gamma}_i + \varepsilon_{ij}$$

Example: Random intercept slope model

$$y_{ij} = [1 \ x_{ij}] \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + [1 \ x_{ij}] \begin{bmatrix} \gamma_{0,i} \\ \gamma_{1,i} \end{bmatrix} + \varepsilon_{ij}$$

Letting

$$\underline{Y}_i = \begin{bmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{in_i} \end{bmatrix}, \quad \underline{X}_i = \begin{bmatrix} \underline{x}_{i1}^\top \\ \underline{x}_{i2}^\top \\ \vdots \\ \underline{x}_{in_i}^\top \end{bmatrix}, \quad \underline{U}_i = \begin{bmatrix} u_{i1}^\top \\ u_{i2}^\top \\ \vdots \\ u_{in_i}^\top \end{bmatrix}, \quad \underline{\varepsilon}_i = \begin{bmatrix} \varepsilon_{i1} \\ \varepsilon_{i2} \\ \vdots \\ \varepsilon_{in_i} \end{bmatrix}$$

then the model is (for cluster  $i=1, 2, \dots, m$ )

$$\underline{Y}_i = \underline{X}_i \underline{\beta} + \underline{U}_i \underline{\gamma}_i + \underline{\varepsilon}_i, \quad \underline{Y}_i \sim N(0, Q)$$

The marginal model is then

$$\underline{Y}_i \sim N\left(\underline{X}_i \underline{\beta}, U_i Q U_i^\top + \sigma^2 I_{n_i}\right)$$

Conditional model

$$\underline{Y}_i | \underline{\gamma}_i \sim N\left(\underline{X}_i \underline{\beta} + U_i \underline{\gamma}_i, \sigma^2 I_{n_i}\right)$$

Letting

$$\underline{Y} = \begin{bmatrix} \underline{y}_1 \\ \underline{y}_2 \\ \vdots \\ \underline{y}_m \end{bmatrix}, \quad \underline{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_m \end{bmatrix}, \quad \underline{\gamma} = \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_m \end{bmatrix}, \quad \underline{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_m \end{bmatrix},$$

$$U = \text{blockdiag}(U_1, U_2, \dots, U_m),$$

the the whole model is

$$\underline{Y} = \underline{X}\beta + U\underline{\gamma} + \underline{\varepsilon}, \quad \underline{\varepsilon} \sim N(0, \underbrace{\sigma^2 I}_{\sum n_i}) = R$$

$$\underline{Y} = \underline{X}\beta + U\underline{\gamma} + \underline{\varepsilon}, \quad \underline{\varepsilon} \sim N(0, \underbrace{\sigma^2 I}_{\sum n_i})$$

where

$$\underline{\gamma} \sim N(0, G), \quad G = \text{blockdiag}(Q_1, Q_2, \dots, Q_m).$$

Conditional model

$$\underline{Y} | \underline{\gamma} \sim N(\underline{X}\beta + U\underline{\gamma}, R)$$

Marginal

$$\underline{Y} \sim N(\underline{X}\beta, \underbrace{U G U^T + R}_{V = V(\underline{\theta})})$$

# Random intercept and slope model, general notation

$$\begin{bmatrix} Y_{11} \\ Y_{12} \\ Y_{13} \\ \vdots \\ Y_{21} \\ Y_{22} \\ \vdots \\ \vdots \\ Y_{m1} \\ Y_{m2} \end{bmatrix} = \begin{bmatrix} 1 & X_{11} \\ 1 & X_{12} \\ \vdots & \vdots \\ 1 & X_{13} \\ X_{21} & \beta_1 \\ X_{22} & \beta_2 \\ \vdots & \vdots \\ \vdots & \vdots \\ 1 & \dots \\ 1 & \end{bmatrix} + \begin{bmatrix} 1 & X_{11} & & & & & & \\ & X_{12} & \cdots & & & & & \\ & & \ddots & & & & & \\ & & & 1 & X_{21} & & & \\ & & & & \ddots & & & \\ & & & & & 1 & X_{22} & \\ & & & & & & \ddots & \\ & & & & & & & \ddots \\ & & & & & & & & 1 & X_{m1} \\ & & & & & & & & & \ddots \\ & & & & & & & & & & 1 & X_{m2} \end{bmatrix} + \begin{bmatrix} \gamma_{0,1} \\ \gamma_{1,1} \\ \gamma_{0,2} \\ \gamma_{1,2} \\ \vdots \\ \vdots \\ \gamma_{0,m} \\ \gamma_{1,m} \end{bmatrix} + \begin{bmatrix} \varepsilon_{11} \\ \varepsilon_{12} \\ \varepsilon_{13} \\ \vdots \\ \varepsilon_{21} \\ \varepsilon_{22} \\ \vdots \\ \vdots \\ \varepsilon_{m1} \\ \varepsilon_{m2} \end{bmatrix}$$

Dimensions:  $n \times 1$     $n \times 2$     $2 \times 1$     $n \times 2m$     $2m \times 1$     $n \times 1$

$$n = \sum_{i=1}^m n_i$$

$$\underline{\gamma} \sim N(0, G), \quad G = \begin{bmatrix} \tau_0^2 & \tau_{01}^2 & & & \\ \tau_{01}^2 & \tau_1^2 & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & \ddots \end{bmatrix}, \quad \varepsilon \sim N(0, R), \quad R = \sigma^2 I_n$$

$$\underline{\theta} = (\tau^2, \tau_0^2, \tau_1^2, \tau_{01}^2)^T$$

## General linear model (Ch. 4, 1)

↑

21.10

$$\underline{y} = \underline{X}\underline{\beta} + \underline{\varepsilon}, \quad \underline{\varepsilon} \sim N(\underline{0}, \sigma^2 \underline{W}^{-1})$$

Multiplying by  $\underline{W}^{1/2}$

$$\underline{W}^{1/2} \underline{y} = \underline{W}^{1/2} \underline{X}\underline{\beta} + \underline{W}^{1/2} \underline{\varepsilon}$$

$$\underline{y}^* = \underline{X}^* \underline{\beta} + \underline{\varepsilon}^*$$

where

$$\text{Var}(\underline{\varepsilon}^*) = \text{Var}(\underline{W}^{1/2} \sigma \underline{W}^{-1} \underline{W}^{1/2} \underline{\varepsilon}) = \sigma^2 \underline{I}_n$$

Hence the MLE of  $\underline{\beta}$  is

$$\begin{aligned}\hat{\underline{\beta}}_{\text{GLS}} &= (\underline{X}^{*T} \underline{X}^*)^{-1} \underline{X}^{*T} \underline{y}^* \\ &= (\underline{X}^T \underline{W} \underline{X})^{-1} \underline{X}^T \underline{W} \underline{y}\end{aligned}$$

$\text{Var}(\hat{\underline{\beta}}_{\text{OLS}})$  is less efficient than  $\text{Var}(\hat{\underline{\beta}}_{\text{GLS}})$  in the sense that  
 $\text{Var}(\hat{\underline{\beta}}_{\text{OLS}}) - \text{Var}(\hat{\underline{\beta}}_{\text{GLS}})$  is positive semi-definite (exercise 21).

## 7.3. ML and REML inference

For  $V = V(\underline{\theta})$  and  $\underline{\theta}$  known the MLE of  $\underline{\beta}$  is

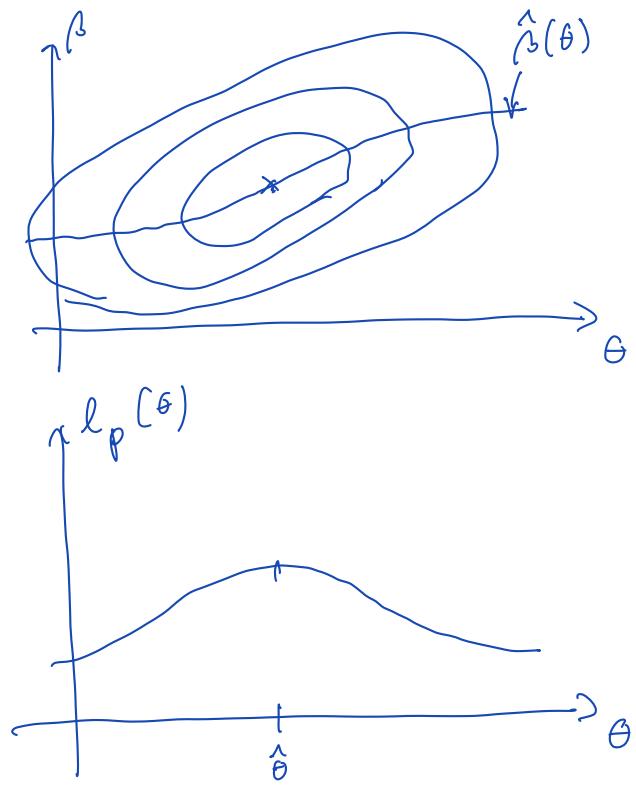
$$\hat{\underline{\beta}} = (\underline{X}^T V^{-1}(\underline{\theta}) \underline{X})^{-1} \underline{X}^T V^{-1}(\underline{\theta}) \underline{y}$$

The profile likelihood for  $\underline{\theta}$  is

$$l_p(\underline{\theta}) = \max_{\underline{\beta}} l(\underline{\theta}, \underline{\beta}) = l(\underline{\theta}, \hat{\underline{\beta}}(\underline{\theta}))$$

$$= -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \ln |V(\underline{\theta})| + \frac{1}{2} (\underline{y} - \underline{X} \hat{\underline{\beta}}(\underline{\theta}))^T V^{-1}(\underline{\theta}) (\underline{y} - \underline{X} \hat{\underline{\beta}}(\underline{\theta}))$$

Maximising this numerically (e.g. using optim in R)  
we obtain the joint MLE of  $\underline{\beta}, \underline{\theta}$ ,  $\hat{\underline{\beta}}(\underline{\theta}), \hat{\underline{\theta}}$



But  $\hat{\theta}_{MLE}$  is typically not unbiased.

Restricted/residual/reduced maximum likelihood (REML)  
(Patterson & Thompson 1971)

Motivating example: Suppose  $y_1, y_2, \dots, y_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

MLE of  $\sigma^2$  if  $\mu$  is known,  $\hat{\sigma}^2 = \frac{1}{n} \sum (y_i - \mu)^2$ ,

is unbiased. But if estimating  $\mu$  by  $\hat{\mu} = \bar{y}$ ,

the joint MLE of  $\sigma^2$ ,  $\hat{\sigma}^2 = \frac{1}{n} \sum (y_i - \bar{y})^2$

is not unbiased. (squared deviation from  $\bar{y}$  smaller than from  $\mu$ )

Idea: Instead of the likelihood of  $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$   
we will work with the likelihood of

$n-p$  "error contrasts"  $\underline{w} = \mathbf{A}^T \underline{y}$  defined such

that  $\mathbf{A}$  is  $n \times (n-p)$ ,  $\mathbf{A}^T \mathbf{A} = \mathbf{I}_{n-p}$

$$\mathbf{S} = \mathbf{A} \mathbf{A}^T = \mathbf{I} - \mathbf{H} = \mathbf{I} - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

The reduced data

$$\underline{w} = \bar{A}^T \underline{y} = \bar{A}^T (\underline{x}\beta + \underline{v}\gamma + \underline{\varepsilon})$$

$$= \cancel{\underline{A}^T \beta} + \bar{A}^T (\underline{v}\gamma + \underline{\varepsilon})$$

$$\sim N(0, \underbrace{\bar{A}^T V(\theta) A}_{\text{becomes dense and expensive to invert in general?}})$$

i.e. the distribution of  $\underline{w}$  doesn't depend on  $\beta$ .

Example cont.,  $\underline{y} \sim N(\mu, \sigma^2 I_n)$ . We then have

$$\text{Var}(\underline{w}) = \text{Var}(\bar{A}^T \underline{y}) = \sigma^2 \bar{A}^T A = \sigma^2 I_{n-p}$$

$$E(\underline{w}) = E(\bar{A}^T \underline{y}) = E(\bar{A}^T (\underline{x}\beta + \underline{v}\gamma + \underline{\varepsilon})) = 0$$

Restricted log likelihood

$$l_R(\sigma^2) = -\frac{n-1}{2} \ln(2\pi) - \frac{1}{2} \ln \left| \sigma^2 I_{n-p} \right| - \frac{1}{2\sigma^2} (\bar{A}^T \underline{y})^T \bar{A}^T \underline{y}$$

$$= -\frac{n-1}{2} \ln(2\pi) - \frac{n-1}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \underline{y}^T (A A^T) \underline{y}$$

$$= \underbrace{\dots}_{\sim \sigma^2} - \underbrace{\frac{1}{2\sigma^2} \underline{y}^T (I - H) \underline{y}}$$

$$\underbrace{\left[ (I - H) \underline{y} \right]^T}_{\hat{\varepsilon}^T \hat{\varepsilon}} (I - H) \underline{y}$$

$$= SSE$$

$$\frac{\partial l_R}{\partial \sigma^2} = 0 \Rightarrow \hat{\sigma}^2_{REML} = \dots = \frac{1}{n-1} SSE = \frac{1}{n-1} \hat{\varepsilon}^T \hat{\varepsilon}$$

$$= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

How to find A?

Let  $S = I - H = I - X(X^T X)^{-1} X^T$ . Eigenvalues and vectors of  $H$  satisfies

$$\lambda \underline{v} = S \underline{v}$$

$$\lambda^2 \underline{v} = S^2 \underline{v} = S \underline{v}$$

since  $S^2 = S$  (why?). Hence

$$S \underline{v} = \lambda \underline{v} = \lambda^2 \underline{v}$$

and all eigenvalues satisfies  $\lambda = \lambda^2$  and are thus either 0 or 1.

For eigenvalues  $\lambda_i = 1$  ( $i = 1, \dots, n-p$ ) the eigenvector satisfies

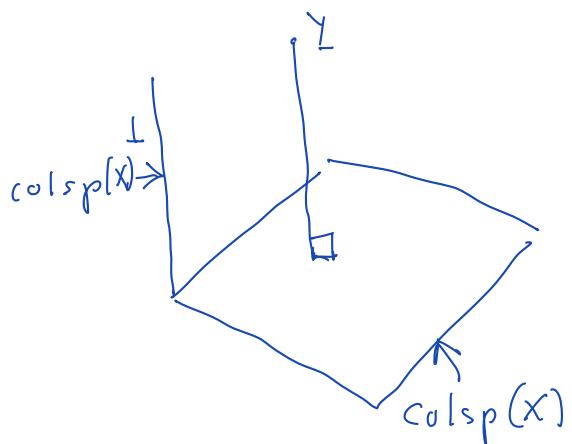
$$S \underline{v}_i = 1 \cdot \underline{v}_i = \underline{v}_i,$$

i.e.  $\underline{v}_i$  is in  $\text{colsp}(X)^\perp$

For eigenvalues  $\lambda_i = 0$  ( $i = n-p+1, \dots, n$ )

$$S \underline{v}_i = 0$$

i.e.  $\underline{v}_i$  is in  $\text{colsp}(X)$ .



Thus, letting

$$\underline{w} = \underline{A}^T \underline{y} = \begin{bmatrix} v_1 & v_2 & \cdots & v_{n-p} \end{bmatrix}^T \underline{y}$$

new  
coordinates  
 $(n-p) \times 1$

new basis vectors

old coordinates  
 $n \times 1$

we have  $\underline{A}^T \underline{x} = 0$ ,  $\underline{A}^T \underline{A} = \mathbb{I}_{n-p}$  ( $S$  is symmetric) and

$$\underline{A} \underline{A}^T = \underline{P} \underline{D} \underline{P}^T$$

$$= \underline{P} \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 0 \\ & & & & 0 \\ & & & & & 0 \end{bmatrix} \underline{P}^T = S = \mathbb{I} - H$$

eigen decomposition  
of  $S$

R:

`A <- eigen(I - X(X^T X)^{-1} X)$vectors[, 1:(n-p)]`

"Bayesian" interpretation / definition of REML (Harville 1974)

We shall see that the restricted likelihood can be expressed as

$$L_R(\underline{\theta}) = \int L(\underline{\beta}, \underline{\theta}) d\underline{\beta}$$

def. in  
(Fahrmeir )

$$\begin{array}{l} \hat{\underline{\beta}} = G^T \underline{y} \\ \text{p} \times 1 \\ \underline{w} = A^T \underline{y} \\ (n-p) \times 1 \end{array}$$

$$\begin{bmatrix} \underline{w} \\ \hat{\underline{\beta}} \end{bmatrix} = [A \quad G]^T \underline{y}$$

$$G^T = (X^T V^{-1}(\theta) X)^{-1} X^T V^{-1}(\theta)$$

$\hat{\underline{\beta}}$  and  $\underline{w}$  are independent because

$$\begin{aligned} \text{Cov}(\hat{\underline{\beta}}, \underline{w}) &= \text{Cov}(G^T \underline{y}, A^T \underline{y}) \\ &= G^T V(\theta) A \\ &= (X^T V^{-1}(\theta) X)^{-1} \underbrace{X^T V^{-1}(\theta) V(\theta) A}_{A^T X = 0} = 0 \end{aligned}$$

Now

$$L_R(\theta) = f_{\underline{w}}(A^T \underline{y} | \theta) \quad = 1 \text{ regardless of whether we integrate w.r.t. } \underline{\beta} \text{ or } \hat{\underline{\beta}}, \text{ for example, note that } \int \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = 1$$

$$= f_{\underline{w}}(A^T \underline{y} | \theta) \int f_{\hat{\underline{\beta}}}(\hat{\underline{\beta}} | \underline{\beta}, \theta) d\underline{\beta}$$

$$= \int f_{\underline{w}, \hat{\underline{\beta}}}(\hat{\underline{\beta}} | A^T \underline{y}, G^T \underline{y}) d\underline{\beta} \quad \left( \text{since } \underline{w} \text{ and } \hat{\underline{\beta}} \text{ are indep.} \right)$$

$$= \int_A G \left( \int f_{\underline{y}}(\underline{y} | \underline{\beta}, \theta) d\underline{\beta} \right) d\underline{w} \quad \left( \text{or } \hat{\underline{\beta}} \text{ linear transf.} \right)$$

$$= \{A, G\} \int L(\beta, \theta) d\beta$$

does not depend on  $\theta$

= alternative def of  $L_R(\theta)$

"Bayesian" interpretation:

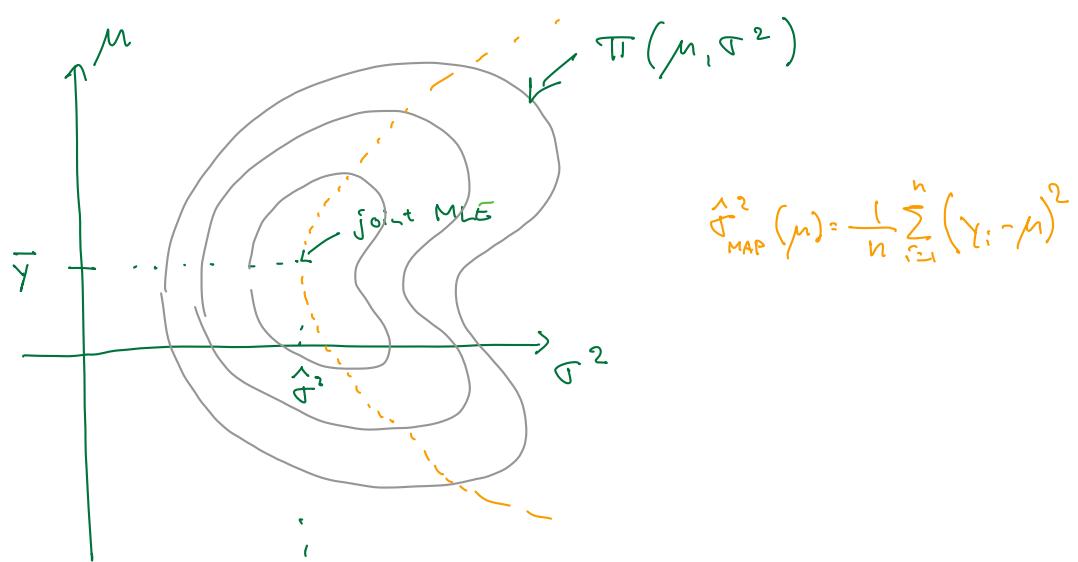
$L(\beta, \theta) \propto$  joint posterior of  $\beta$  and  $\theta$   
if using uniform improper priors on  $\beta$  and  $\theta$

$\int L(\beta, \theta) d\beta \propto$  marginal posterior distribution of  $\theta$

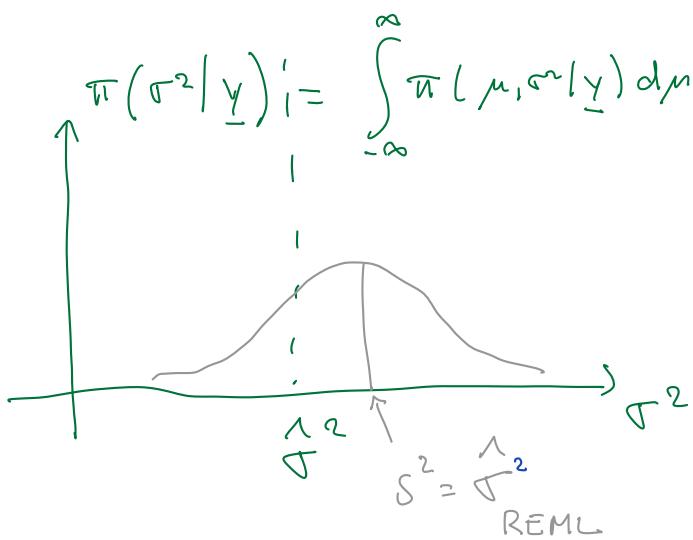
$\hat{\theta}_{REML}$  = Maximum a posteriori (MAP) estimate of  $\theta$ .

Example:  $y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$ .  
Joint posterior of  $\mu, \sigma^2$ :  $= 1$

$$\begin{aligned} \pi(\mu, \sigma^2 | \bar{y}) &\propto f(\bar{y} | \mu, \sigma^2) \pi(\mu, \sigma^2) \\ &= \dots \propto \frac{1}{\sigma^2} e^{-\frac{(\mu - \bar{y})^2}{2\sigma^2/n}} \end{aligned}$$



Marginal posterior of  $\sigma^2$



Why doesn't  $|A G|$  depend on  $\theta$ ?

$$|A G| = \left| \begin{bmatrix} A^\top \\ G^\top \end{bmatrix} \begin{bmatrix} A & G \end{bmatrix} \right|^{1/2}$$

wikipedia

$$\begin{vmatrix} A & B \\ C & D \end{vmatrix} = |A| \times |D - CA^{-1}B|$$

$$A^\top A = I_{n-p}$$

$$\begin{aligned} AA^\top &= I - H \\ &= I - X(X^\top X)^{-1}X^\top \end{aligned}$$

$$= \begin{vmatrix} A^\top A & A^\top G \\ G^\top A & G^\top G \end{vmatrix}^{1/2}$$

$$= |A^\top A|^{1/2} \times |G^\top G - G^\top A (A^\top A)^{-1} A^\top G|^{1/2}$$

$$= |I|^{1/2} \times |G^\top (I - \underbrace{(A^\top A)^{-1}}_{\text{yellow box}}) G|^{1/2}$$

$$= |G^\top (I - X(X^\top X)^{-1}X^\top) G|^{1/2}$$

$$= |(X^\top X)^{-1}|^{1/2} \times |X(X^\top X)^{-1}X^\top|^{1/2} \times |(X^\top X)^{-1}|^{1/2}$$

$$= |(X^\top X)|^{-1/2} \times |X^\top X|^{-1/2}$$

i.e. doesn't depend on  $\theta$  (although  $G$  depends on  $\theta$ )

21-10  
26-10

Connection between profile and restricted likelihood  
 $L_p(\underline{\theta})$  and  $L_R(\underline{\theta})$  (Eq. 7.29 in Fahr)

Recall the profile likelihood

$$L_p(\underline{\theta}) = (2\pi)^{-\frac{n}{2}} |V(\underline{\theta})|^{-\frac{1}{2}} e^{-\frac{1}{2}(\underline{y} - \underline{x}\hat{\beta}(\underline{\theta}))^T V^{-1}(\underline{\theta})(\underline{y} - \underline{x}\hat{\beta}(\underline{\theta}))} \quad (1)$$

The restricted likelihood

$$L_R(\underline{\theta}) = \int f(\underline{y} | \underline{\beta}, \underline{\theta}) d\underline{\beta}$$

$$= (2\pi)^{-\frac{n}{2}} |V(\underline{\theta})|^{-\frac{1}{2}} \int e^{-\frac{1}{2}(\underline{y} - \underline{x}\underline{\beta})^T V^{-1}(\underline{\theta})(\underline{y} - \underline{x}\underline{\beta})} d\underline{\beta} \quad (2)$$

$\underbrace{p \text{ dimension}}$

Consider the decomposition

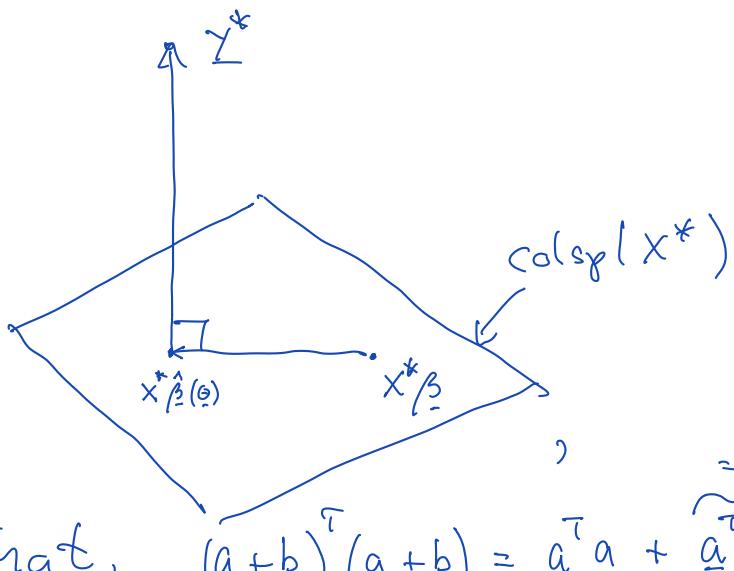
$$V(\underline{\theta})^{-1/2}(\underline{y} - \underline{x}\underline{\beta}) = \underline{y}^* - \underline{x}^*\underline{\beta}$$

$$= \underline{y}^* - \underline{x}^*\hat{\beta}(\underline{\theta}) + \underline{x}^*\hat{\beta}(\underline{\theta}) - \underline{x}^*\underline{\beta}$$

$$= \underline{V}(\underline{\theta})^{-1/2}(\underline{y} - \underline{x}\hat{\beta}(\underline{\theta})) + \underline{V}(\underline{\theta})^{1/2} \underline{x}(\hat{\beta}(\underline{\theta}) - \underline{\beta}) \quad (3)$$

$$= \underline{a} + \underline{b}$$

The vectors on the right hand side are orthogonal:



such that,  $(\underline{a} + \underline{b})^\top (\underline{a} + \underline{b}) = \underline{a}^\top \underline{a} + \underbrace{\underline{a}^\top \underline{b}}_{=0} + \underbrace{\underline{b}^\top \underline{a}}_{=0} + \underline{b}^\top \underline{b}$ .

Hence, inserting (3) into (2) we obtain

$$L_R(\underline{\theta}) = (2\pi)^{-\frac{n}{2}} |V(\underline{\theta})|^{-\frac{1}{2}} e^{-\frac{1}{2}(y - X\hat{\beta}(\underline{\theta}))^\top V^{-1}(\underline{\theta})(y - X\hat{\beta}(\underline{\theta}))}$$

$$\times \int e^{-\frac{1}{2}((\hat{\beta}(\underline{\theta}) - \beta)^\top X^\top V^{-1}(\underline{\theta}) X (\hat{\beta}(\underline{\theta}) - \beta))} d\beta$$

$$= L_p(\underline{\theta}) (2\pi)^{\frac{p}{2}} \left| \underbrace{X^\top V^{-1}(\underline{\theta}) X}_{n \times n} \right|^{-\frac{1}{2}}$$

$$\sim \underbrace{C_1}_{p \times p} \underbrace{\text{normalising constant}}$$

Computing this only involves computing the inverse and determinant of a block diagonal  $n \times n$  matrix which is cheap.

## Hypothesis testing for LMMs

$$\text{Var}(\hat{\beta}) = V(\hat{\theta})$$

We have

$$\hat{\beta} = (X^T V^{-1}(\hat{\theta}) X)^{-1} X^T V^{-1}(\hat{\theta}) \hat{y}$$

Ignoring uncertainty in  $\hat{\theta}$ ,

$$\text{Var}(\hat{\beta}) = \dots = X^T V^{-1}(\hat{\theta}) X$$

is an estimate of  $\text{Var}(\hat{\beta})$ .

The Wald statistic

$$(C\hat{\beta} - d)^T \left( C X^T V^{-1}(\hat{\theta}) X C^T \right)^{-1} (C\hat{\beta} - d)$$

is approximately chi-square. No exact F-test

LRTs

Fixed effect testing: LRTs must be based on the full likelihood because the restricted maximum likelihoods are likelihoods of different subsets of the data.

Random effects: LRT and RLRT based on  $L_R(\hat{\theta})$  can both be used as the likelihoods are based on the same  $\mathbf{w} = A^T \hat{y}$

R: lmer

$H_0$  on the boundary

Example 1 (random intercept model)

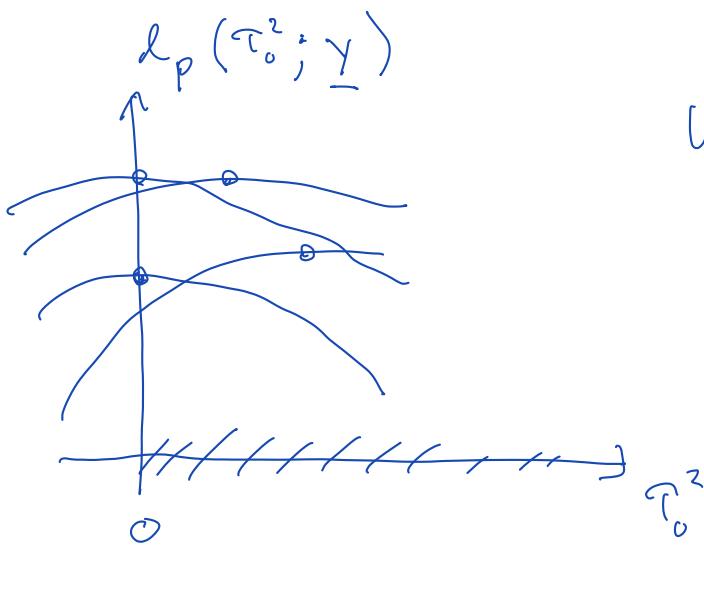
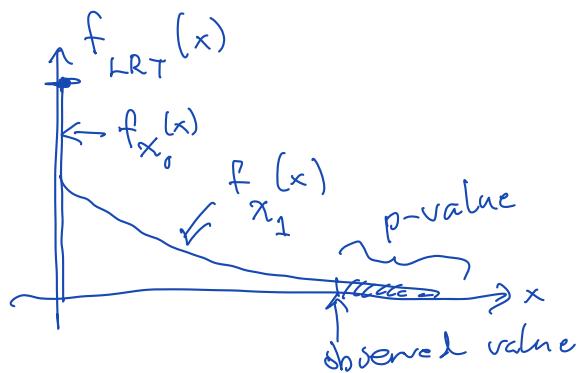
$H_0: \tau_0^2 = 0$  (no random intercept in model)

$H_1: \tau_0^2 > 0$  (random effect present)

$$LRT \approx 2 \left( l_p(\hat{\tau}_0^2) - l_p(0) \right) \stackrel{\text{asym}}{\sim} 0.5 \chi_0^2 : 0.5 \chi_1^2$$

i.e.

$$f_{LRT}(x) = 0.5 \underbrace{\delta(x)}_{\substack{\text{point mass} \\ \text{at zero}}} + 0.5 f_{\chi_1^2}(x)$$



Under  $H_0$

$$P(\hat{\tau}_0^2 > 0) = \frac{1}{2}$$

MLE under  $H_1$

$$\text{since } E\left(\frac{\partial l}{\partial \tau_0^2} \Big| \tau_0^2 = 0\right) = 0$$

## Example 2: Random intercept and slope model

$$H_0: \tau_0^2 \geq 0, \tau_1^2 = 0, \tau_{01}^2 = 0 \quad (\text{only random intercept})$$

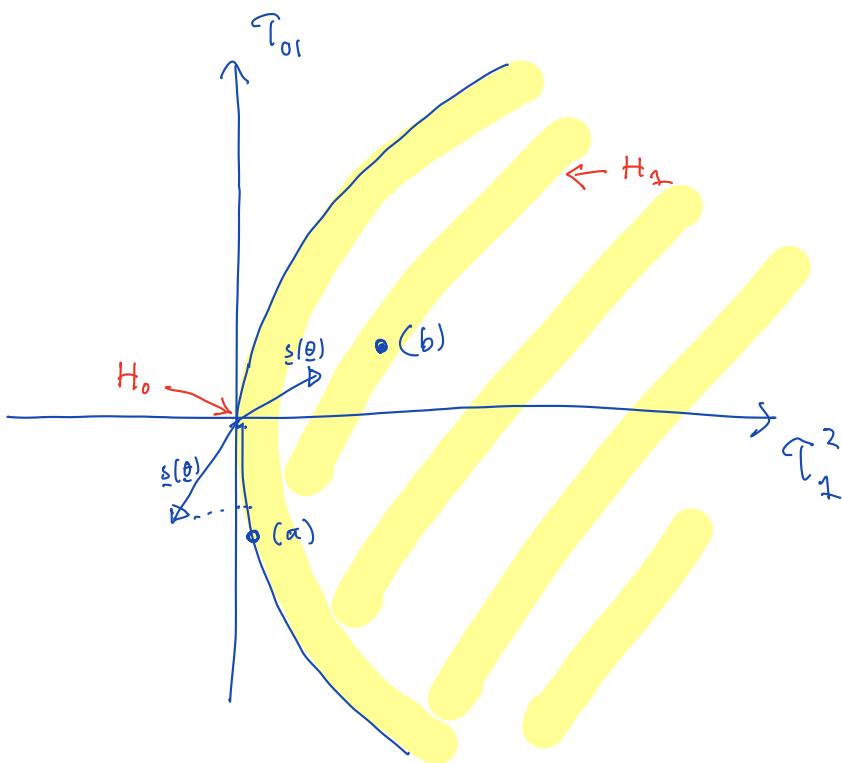
$$H_1: \tau_0^2 \geq 0, \tau_1^2 > 0, |\tau_{01}| < \sqrt{\tau_0^2 \tau_1^2} \quad (\text{random intercept and slope})$$

$$\text{LRT} = 2(\ell(\hat{\tau}_0^2, \hat{\tau}_1^2, \hat{\tau}_{01}^2) - \ell(\hat{\tau}_{0,0}^2, 0, 0))$$

asympt.

$$\sim 0.5 \chi_1^2 : 0.5 \chi_2^2 \quad (\text{Fahrmeir})$$

Heuristic explanation?



Asymptotically, under  $H_0$ , there is a 50% chance that the MLE under  $H_1$  ends up on the boundary (a) and 50% chance that it ends up at the interior (b)

26.10



# Estimating the random effects $\gamma$ (7.3.1, 7.3.3)

↑

Recall conditional model

28.10

$$\underline{y} | \underline{\gamma} \sim N(\underline{x}\underline{\beta} + \underline{v}_j, R)$$

and

$$\underline{\gamma} \sim N(0, G)$$

Given  $\underline{\beta}$  estimated by REML or ML,  $\underline{\alpha} = \begin{bmatrix} \underline{\beta} \\ \underline{\gamma} \end{bmatrix}$  can be estimated by maximizing the "joint penalized likelihood"

$$L(\underline{\alpha}) = L(\underline{\beta}, \underline{\gamma})$$

$$= f_{\underline{y}|\underline{\gamma}}(\underline{y}) f_{\underline{\gamma}}(\underline{\gamma})$$
  
$$\propto e^{-\frac{1}{2} (\underline{y} - \underline{x}\underline{\beta} - \underline{v}_j)^T R^{-1} (\underline{y} - \underline{x}\underline{\beta} - \underline{v}_j) - \frac{1}{2} \underline{\gamma}^T G^{-1} \underline{\gamma}}$$

} Gaussian in  $\underline{\alpha}$

or minimising

$$(\underline{y} - \underline{x}\underline{\beta} - \underline{v}_j)^T R^{-1} (\underline{y} - \underline{x}\underline{\beta} - \underline{v}_j) + \underline{\gamma}^T G^{-1} \underline{\gamma}$$
$$= (\underline{y} - C\underline{\alpha})^T R^{-1} (\underline{y} - C\underline{\alpha}) + \underline{\alpha}^T B \underline{\alpha} \quad = (*)$$

where

$$C = \underbrace{[\underline{x} \ \underline{v}_j]}_{n \times (p+1)} \quad \text{and} \quad B = \begin{bmatrix} 0 & 0 \\ 0 & G^{-1} \end{bmatrix}$$

Recall

$$\frac{\partial}{\partial \underline{x}} (\underline{x}^T A \underline{x}) = 2A\underline{x}$$

$$\frac{\partial}{\partial \underline{x}} (A \underline{x})^T = A^T$$

$$\frac{\partial}{\partial \underline{x}} f(g(\underline{\alpha})) = \frac{\partial g^T}{\partial \underline{x}} \cdot \frac{\partial f}{\partial g}$$

Thus

$$\frac{\partial L}{\partial \underline{\alpha}} = -C^T 2R(\underline{y} - C\underline{\alpha}) + 2B\underline{\alpha}$$
$$= -2C^T R^{-1}\underline{y} + 2(C^T R^{-1}C + B)\underline{\alpha}$$

Equating to 0 and solving

$$\hat{\underline{\alpha}} = (C^T R^{-1}C + B)^{-1} C^T R^{-1} \underline{y}$$

Since  $L(\underline{\alpha}) \propto$  Gaussian function of  $\underline{\alpha}$  with variance matrix  $(C^T R^{-1}C + B)^{-1}$ . Thus, the "posterior variance" of  $\underline{\alpha}$

$$\text{Var}(\underline{\alpha} | \underline{y}) = (C^T R^{-1}C + B)^{-1}$$

If instead fixing  $\underline{x}$  and thinking of  $\underline{y}$  as random (frequentist view)

$$\begin{aligned} \text{Var}(\hat{\underline{\alpha}}) &= \text{Var}((C^T R^{-1}C + B)^{-1} C^T R^{-1} \underline{y}) \\ &= (C^T R^{-1}C + B)^{-1} C^T R^{-1} R R^{-1} C (C^T R^{-1}C + B)^{-1} \\ &= \text{Var}(\underline{\alpha} | \underline{y}) C^T R^{-1} C (C^T R^{-1}C + B)^{-1} \\ &\leq \text{Var}(\underline{\alpha} | \underline{y}), \text{ i.e. } \text{Var}(\underline{\alpha} | \underline{y}) - \text{Var}(\underline{\alpha}) \text{ is p.s.d.} \end{aligned}$$

Alternative formula: Recall that if

$$\begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_n \end{bmatrix} \sim N\left(\begin{bmatrix} \underline{\mu}_1 \\ \vdots \\ \underline{\mu}_n \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \vdots & \vdots \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

then

$$E(\underline{x}_1 | \underline{x}_2) = \underline{\mu}_1 + \Sigma_{12} \Sigma_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2)$$

We have

$$\begin{bmatrix} \underline{y} \\ \vdots \\ \underline{y} \end{bmatrix} \sim N\left(\begin{bmatrix} \underline{0} \\ \vdots \\ \underline{X}\beta \end{bmatrix}, \begin{bmatrix} G & G U^T \\ U G & \underbrace{U G U^T + R}_{=V} \end{bmatrix}\right)$$

$\text{cov}(A\underline{x}, B\underline{y})$

$$Y = X\beta + U\underline{x} + \underline{\varepsilon}$$

Hence

$$\hat{\underline{y}} = E(\underline{y} | \underline{y}) = \underline{0} + G U^T \hat{V}^{-1} (\underline{y} - X \hat{\beta}) \quad (\text{BLUP})$$

Example: Random intercept model: Observation  $j=1, \dots, n_i$  in cluster  $i$

$$Y_{ij} = \underline{x}_{ij}^T \hat{\beta} + \gamma_{0,i} + \varepsilon_{ij}$$

$$\underline{Y}_i = \underline{X}_i^T \hat{\beta} + \underbrace{\underline{U}_i \gamma_{0,i}}_{\frac{1}{n_i}} + \underline{\varepsilon}_i, \quad \gamma_{0,i} \sim N(0, \tau_0^2), \quad \underline{\varepsilon}_i \sim N(\underline{0}, \sigma^2 I_{n_i})$$

$$\begin{aligned} V_i &= \text{Var}(\underline{y}_i) = \underline{U}_i \tau_0^2 \underline{U}_i^T + \sigma^2 I_{n_i} \\ &= \tau_0^2 \underbrace{J}_{n_i \times n_i} + \sigma^2 I_{n_i} \end{aligned}$$

$$\hat{A}^{-1} A = \mathbb{I}$$

We need

$$V_i^{-1} = a I_{n_i} + b J_{n_i}$$

satisfy 2)

$$(\sigma^2 I + \tau_0^2 J)(a I + b J) = \mathbb{I}$$

$$J = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \ddots \end{bmatrix}$$

$$J^2 = J J = ?$$

$$a \sigma^2 I + \underbrace{\sigma^2 b J}_{\sigma^2 b n_i J} + \underbrace{\tau_0^2 a J}_{\tau_0^2 a n_i J} + \underbrace{\tau_0^2 b n_i J}_{\tau_0^2 b n_i J} = \mathbb{I}$$

$$a \sigma^2 = 1 \Rightarrow a = \frac{1}{\sigma^2}$$

$$\Rightarrow b = - \frac{\tau_0^2}{(\sigma^2 + n_i \tau_0^2) \sigma^2}$$

$$V_i^{-1} = \frac{1}{\sigma^2} \left( I_{n_i} - \frac{\tau_0^2}{\sigma^2 + n_i \tau_0^2} J_{n_i} \right)$$

We have

$$\hat{f} = G U^\top \hat{V}^{-1} (\underline{y} - \underline{x}\hat{\beta})$$

and within each cluster  $i = 1, \dots, m$

$$\hat{f}_i = \tau_0^2 \frac{1^\top}{n_i} \frac{1}{\sigma^2} \left( I_{n_i} - \frac{\tau_0^2}{\sigma^2 + n_i \tau_0^2} J_{n_i} \right) (\underline{y}_i - \underline{x}_i \hat{\beta})$$

$$= \frac{\tau_0^2}{\sigma^2} \left( \frac{1^\top}{n_i} - \frac{1}{\sigma^2 + n_i \tau_0^2} \right) \underline{y}_i - \underline{x}_i \hat{\beta}$$

$$= \frac{\tau_0^2}{\sigma^2} \left( 1 - \frac{\tau_0^2 n_i}{\sigma^2 + n_i \tau_0^2} \right) n_i \frac{1}{n_i} \sum_{j=1}^{n_i} (y_{ij} - x_{ij} \hat{\beta})$$

$$= \underbrace{\frac{n_i \tau_0^2}{\sigma^2 + n_i \tau_0^2}}_{\text{shrinkage factor } < 1} \cdot \underbrace{\frac{1}{n_i} \sum_{j=1}^{n_i} (y_{ij} - x_{ij} \hat{\beta})}_{\text{mean residual within cluster } i}$$

Small  $n_i$  or  $\tau_0^2$   
⇒ more shrinkage

R: mod ← lmer (Reaction ~ 1 + Days + (1 + Days | subject))

gammahat ← ranef(mod, condVar = TRUE)

gammahat\$subject #  $\hat{f}$

attributed(gammahat\$subject) # Var( $\hat{f} | y$ )  
as a  $m \times 2 \times 2$  array

Uncertainty in  $\underline{\theta} = (\sigma^2, \tau_0^2, \tau_i^2, \tau_{0i}^2)$  ignored.

## Generalized linear mixed model (GLMMs) (Ch. 7.5)

Distribution of the response  $y_{ij}$  conditional on random effect  $\gamma_i$ ,  $f(y_{ij} | \gamma_i)$  belongs to the exponential family with

$$\mu_{ij} = E(y_{ij} | \gamma_i)$$

related to

$$\eta_{ij} = \underline{x}_{ij}^\top \underline{\beta} + u_{ij}^\top \underline{\gamma}_i, \quad \underline{\gamma}_i \sim N(\underline{0}, Q)$$

for  $i=1, 2, \dots, m$  and  $j=1, 2, \dots, n_i$  through

$$\mu_{ij} = h(\eta_{ij}) \quad \text{or} \quad g(\mu_{ij}) = \eta_{ij}.$$

All  $y_{ij}$  are conditionally independent.

### Conditional versus marginal models (7.5.2)

Marginal distribution of  $y_{ij}$  becomes

$$f(y_{ij}) = \int f(y_{ij} | \gamma_i) d\gamma_i$$

$$= \int f(y_{ij} | \gamma_i) f(\gamma_i) d\gamma_i$$

that is, a compound distribution or uncountable mixture of e.g. binomial or poisson distributions

logitnormal binomial

lognormal Poisson

Marginal joint distribution of  $\underline{y}_i = (y_{i1}, y_{i2}, \dots, y_{in_i})$

$$f(\underline{y}_i) = \int f(\underline{y}_i | \underline{\gamma}_i) f(\underline{\gamma}_i) d\underline{\gamma}_i$$

$$= \prod_{j=1}^{n_i} f(y_{ij} | \gamma_j) f(\gamma_j) d\gamma_j$$

must be approximated / computed numerically

Poisson random intercept model

$$\begin{aligned} E(y_{ij} | \gamma_{0i}) &= e^{\gamma_{0i}} = e^{x_{ij}^\top \beta + u_{ij}^\top \delta} \\ &= e^{x_{ij}^\top \beta + \gamma_{0i}} = e^{x_{ij}^\top \beta} e^{\gamma_{0i}} \end{aligned}$$

$Z \sim \text{lognormal}(\mu, \sigma^2)$

$$\ln Z = X \sim N(\mu, \sigma^2), M_X(t) = e^{\mu t + \frac{1}{2} \sigma^2 t^2}$$

$$E Z = E(e^X) = M_X(1) = e^{\mu + \frac{1}{2} \sigma^2}$$

$$E(Z^2) = E(e^{2X}) = M_X(2) = e^{2\mu + 2\sigma^2}$$

$$\text{Var } Z = E(Z^2) - (EZ)^2 = \dots = \underbrace{[EZ]^2 (e^{\sigma^2} - 1)}$$

Law of total expectation (LTE)  
 $E(E(X|Y)) = EX$

Meaning of  $EX|Y$ ?

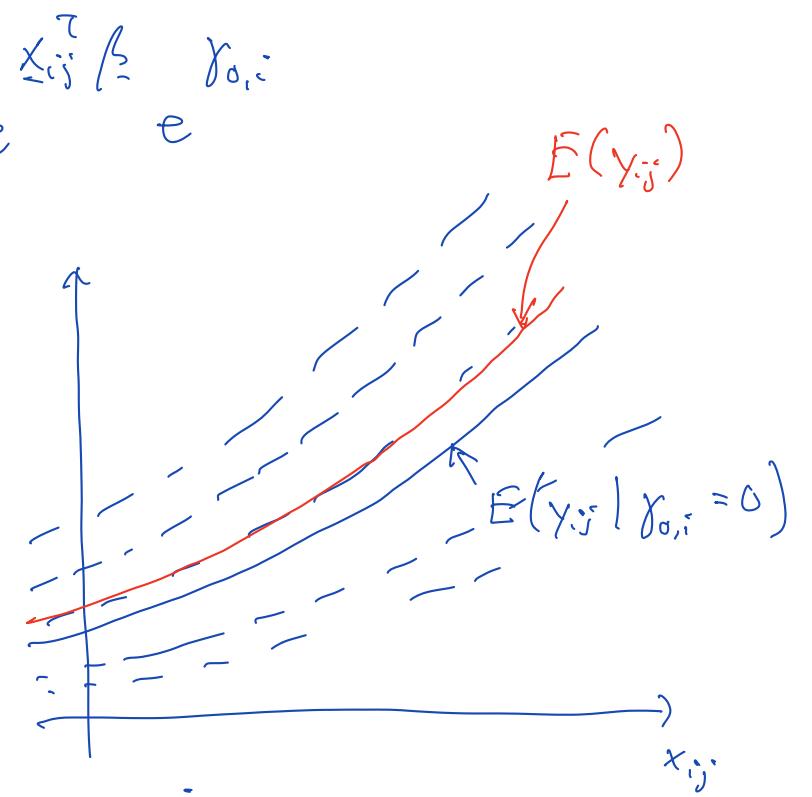
Conditional expectation

$$E(X|Y=x) = g(x)$$

By  $E(X|Y)$  we mean the random  $g(Y)$

Hence, using LTE,

$$\begin{aligned} EY_{ij} &= E(EY_{ij}|Y_{0,i}) \leq Ee^{\beta_0 Y_{0,i}} \\ &= e^{X_{ij}^\top \beta + E[e^{\beta_0 Y_{0,i}}]} \\ &= e^{X_{ij}^\top \beta + \left(1 + \frac{1}{2}\sigma_0^2\right)} \end{aligned}$$



Effect of unit change in  $X_{ijk}$  is that both  $E(Y_{ij})$  and  $E(Y_{ij}|Y_{0,i})$  increase by factor  $e^{\beta_k}$

Using the law of total variance,

$$\begin{aligned} \text{Var} Y_{ij} &= E \text{Var}(Y_{ij}|Y_{0,i}) + \text{Var}(EY_{ij}|Y_{0,i}) \\ &= E e^{X_{ij}^\top \beta + \gamma_{0,i}} + \text{Var} \left( e^{X_{ij}^\top \beta + \gamma_{0,i}} \right) \end{aligned}$$

$$= e^{\frac{x_{ij}^\top \beta}{\sigma_e^2} + \frac{1}{2}\tau_0^2} + e^{\frac{2x_{ik}^\top \beta}{\sigma_e^2}} e^{\tau_0^2} (e^{\sigma_0^2} - 1)$$

Using law of total covariance, the intraclass covariance

$$\begin{aligned} \text{Cov}(y_{ij}, y_{ik}) &= E \text{Cov}(y_{ij}, y_{ik} | \gamma_{0,i}) + \text{Cov}[E(y_{ij} | \gamma_{0,i}), E(y_{ik} | \gamma_{0,i})] \\ &= E 0 + \text{Cov}\left(e^{\frac{x_{ij}^\top \beta}{\sigma_e^2} + \gamma_{0,ij}}, e^{\frac{x_{ik}^\top \beta}{\sigma_e^2} + \gamma_{0,ik}}\right) \\ &= e^{(x_{ij}^\top + x_{ik}^\top)\beta / \sigma_e^2} \text{Var}(e^{\gamma_{0,ij}}) \\ &= e^{(x_{ij}^\top + x_{ik}^\top)\beta / \sigma_e^2} e^{\tau_0^2} (e^{\sigma_0^2} - 1). \end{aligned}$$

Intraclass correlation:

$$\text{corr}(y_{ij}, y_{ik}) = \dots$$

### Binary random intercept probit model

Conditional model

$$\text{probit } P(y_{ij} = 1 | \gamma_{0,i}) = x_{ij}^\top \beta + \gamma_{0,ij}$$

or

$$P(y_{ij} = 1 | \gamma_{0,i}) = \Phi(x_{ij}^\top \beta + \gamma_{0,ij}), \quad \gamma_{0,ij} \sim N(0, \tau_0^2)$$

Marginal model

$$P(y_{ij} = 1) \stackrel{\text{LTP}}{=} \int P(y_{ij} = 1 | \gamma_{0,i}) f(\gamma_{0,i}) d\gamma_{0,i}$$

$$= \int \Phi(x_{ij}^\top \beta + \gamma_{0,ij}) f(\gamma_{0,i}) d\gamma_{0,i}$$

$$= \int P(Z \leq x_{ij}^\top \beta + \gamma_{0,i} \mid \delta_{0,i}) f(\gamma_{0,i}) d\gamma_{0,i}$$

LTP

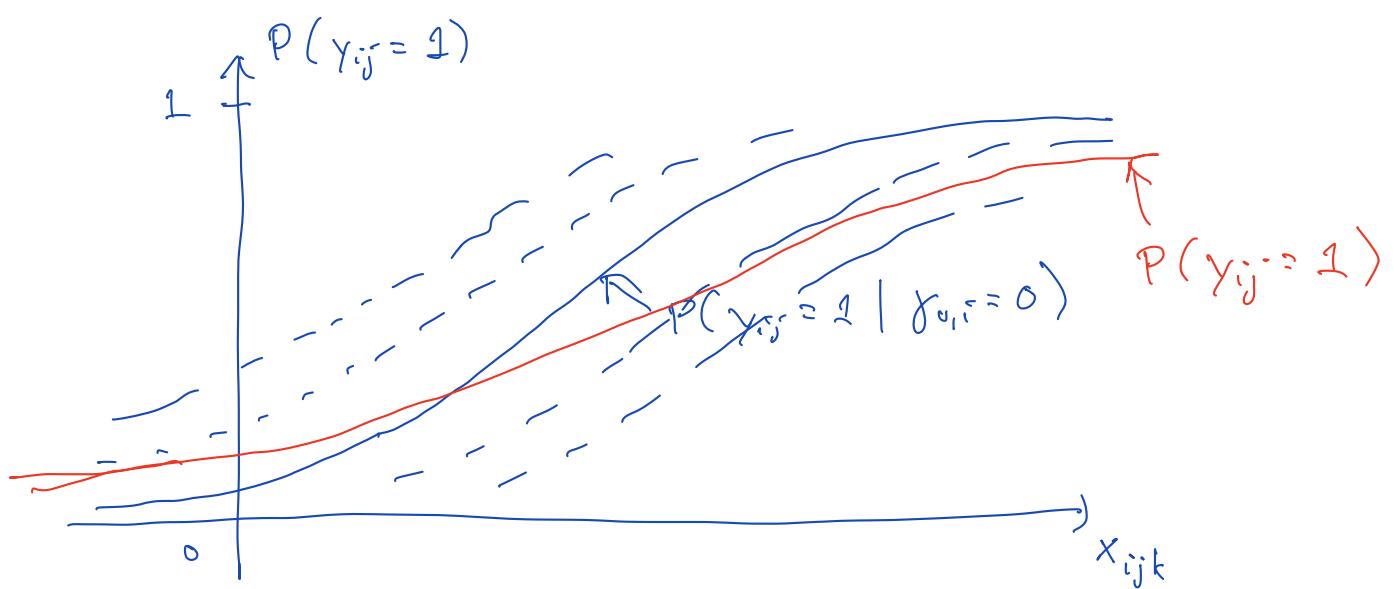
$$\leq P(Z \leq x_{ij}^\top \beta + \gamma_{0,i})$$

$$= P\left(\frac{Z - \gamma_{0,i}}{\sqrt{1 + \tau_0^2}} \leq \frac{x_{ij}^\top \beta}{\sqrt{1 + \tau_0^2}}\right)$$

$\sim N(0, 1)$

$$= \Phi\left(\frac{x_{ij}^\top \beta}{\sqrt{1 + \tau_0^2}}\right)$$

The effect is that  $\beta$  is deflated by  
a factor  $\frac{1}{\sqrt{1 + \tau_0^2}}$



### Binary random intercept logit model

Conditionally

$$P(y_{ij} = 1 | \gamma_{0,i}) = \frac{1}{1 + e^{-(\bar{x}_{ij}^\top \beta + \gamma_{0,i})}}$$

Marginally

$$P(y_{ij} = 1) \approx \frac{1}{1 + e^{-\bar{x}_{ij}^\top \beta / \sqrt{1 + 0.6 \tau^2}}}$$

(Agresti 2002, p. 499)

## Inference for GLMMs

1) Numerical integration      `lme4::glmer( , nAGP = 10 )`

$$L(\beta, \theta) = f(\underline{y} | \beta, \theta)$$

$$= \prod_{i=1}^m f(\underline{y}_i | \beta, \theta) \quad (\text{clustered data})$$

$$= \prod_{i=1}^m \int f(\underline{y}_i | \underline{\gamma}_i, \beta) f(\underline{\gamma}_i | \theta) d\underline{\gamma}_i \quad (LTP)$$

$$= \prod_{i=1}^m \left( \int \underbrace{\prod_{j=1}^{n_i} f(y_{ij} | \underline{\gamma}_{ij}, \beta)}_{\int_{-\infty}^{\infty} g(t) dt} f(\underline{\gamma}_i | \theta) d\underline{\gamma}_i \right)$$

Numerical integration of  $g(t)$  via Gauss quadrature:

# Laplace approximation example (Stirling's formula)

$$n! = \Gamma(n+1)$$

$$g(x) = x - n \ln x$$

$$= \int_0^\infty x^n e^{-x} dx$$

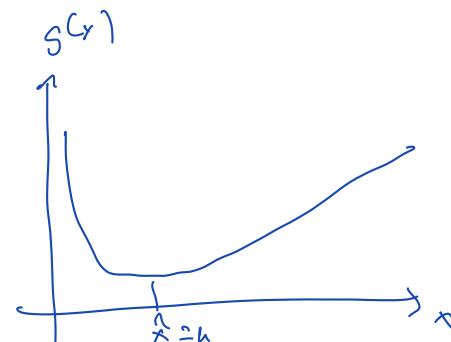
$$g'(x) = 1 - \frac{n}{x} = 0$$

Minimum at  $\hat{x} = n$

$$= \int_0^\infty e^{-(x - n \ln x)} dx$$

$$= \int_0^\infty e^{-g(x)} dx$$

$$= \int_0^\infty e^{-g(n) - \frac{1}{2} g''(n)(x-n)^2} dx$$



$$g''(x) = \frac{n}{x^2} = \frac{1}{n}$$

$$= e^{-n + n \ln n} \int_0^\infty e^{-\frac{1}{2n}(x-n)^2} dx$$

$$\stackrel{n \text{ large}}{=} e^{-n + n \ln n} \int_{-\infty}^\infty e^{-\frac{1}{2} \left( \frac{x-n}{\sqrt{n}} \right)^2} dx$$

$$= e^{-n + n \ln n} \sqrt{2\pi n}$$

$$= \underbrace{\sqrt{2\pi n}}_{\sim \sqrt{n}} e^{-n}$$

2.11



2) Laplace approximation of marginal likelihood

↑  
4.11

$$\begin{aligned}
 L(\beta, \theta) &= \int f(y | \underline{\delta}, \beta) f(\underline{\delta} | \theta) d\underline{\delta} \\
 &= \int e^{-g(\underline{\delta}, \beta, \theta)} f(y, \underline{\delta} | \beta, \theta) d\underline{\delta} \\
 &\approx \int e^{-g(\hat{\delta}(\beta, \theta), \beta, \theta) - \frac{1}{2} (\underline{y} - \hat{\delta}(\beta, \theta))^T H(\beta, \theta) (\underline{y} - \hat{\delta}(\beta, \theta))} d\underline{\delta}
 \end{aligned}$$

where  $\hat{\delta}(\beta, \theta) = \underset{\underline{\delta}}{\operatorname{argmin}} g(\underline{\delta}, \beta, \theta)$  (inner optimization)

and

$$H(\beta, \theta) = \left. \frac{\partial^2}{\partial \underline{\delta} \partial \underline{\delta}^T} g(\underline{\delta}, \beta, \theta) \right|_{\underline{\delta} = \hat{\delta}(\beta, \theta)}$$

Thus,

$$L(\beta, \theta) \approx e^{-g(\hat{\delta}(\beta, \theta), \beta, \theta) - \frac{1}{2} (z\bar{u})^T H(\beta, \theta)^{-1} (z\bar{u})}$$

where  $k$  is the number of random effects

### Automatic differentiation

Used by R-packages TMB and glmmTMB,

Stan (different from symbolic and numerical differentiation)

Example : Find  $\frac{\partial f}{\partial \underline{x}}$  for  $f(\underline{x}) = f(x_1, x_2, x_3)$

$$= \left( x_1 x_2 \sin(x_3) + e^{x_1 x_2} \right) / x_3$$

$$x_1 = 1$$

$$x_2 = 2$$

$$x_3 = \frac{\pi}{2}$$

$$w_1 = x_1 x_2$$

$$w_2 = \sin(x_3)$$

$$w_3 = \exp(w_1)$$

$$w_4 = w_1 w_2$$

$$w_5 = w_3 + w_4$$

$$w_6 = w_5 / x_3$$

|       |         | $\frac{\partial}{\partial x_1}$ | $\frac{\partial}{\partial x_2}$ | $\frac{\partial}{\partial x_3}$ |
|-------|---------|---------------------------------|---------------------------------|---------------------------------|
| $x_1$ | 1       | 1                               | 0                               | 0                               |
| $x_2$ | 2       | 0                               | 1                               | 0                               |
| $x_3$ | $\pi/2$ | 0                               | 0                               | 1                               |
| $w_1$ | 2       | 2                               | 1                               | 0                               |
| $w_2$ | 1       | 0                               | 0                               | 0                               |
| $w_3$ | -       | -                               | -                               | -                               |

dual numbers

$$w_1 = x_1 x_2 = 2$$

$$\frac{\partial w_1}{\partial \underline{x}} = \frac{\partial}{\partial \underline{x}} (x_1 x_2)$$

$$= \frac{\partial x_1}{\partial \underline{x}} x_2 + x_1 \frac{\partial x_2}{\partial \underline{x}}$$

$$= \begin{bmatrix} 1 \\ 0 \end{bmatrix} \cdot 2 + 1 \cdot \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$w_6 \left| \begin{array}{c} (1) \\ (2) \end{array} \right|$$

$$w_3 = e^{w_1}$$

$$\underbrace{\frac{\partial w_3}{\partial x}}_{=} = e^{w_1} \frac{\partial w_1}{\partial x} = e^2 \cdot \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}$$

etc...

## REML for GLMMs (Millar 2011)

For LMMs:

$$\underline{L}_R(\underline{\theta}) = f_{\underline{w}}(\underline{w}), \quad \underline{w} = \underline{A}^T \underline{y}$$

$m$        $n$   
 $m \times 1$        $n \times 1$

but this doesn't generalize to GLMMs.

But the equivalent definition

$$\underline{L}_R(\underline{\theta}) = \int \underline{L}(\underline{\theta}, \underline{\beta}) d\underline{\beta}$$

$$= \int \int f(y|x, \beta) f(x|\theta) dx d\beta$$

does and can be evaluated by Laplace approximation w.r.t. both  $x$  and  $\beta$

`glmmTMB( , REML = TRUE )`

## Gauss Legendre / Hermite quadrature

Warm up: If  $f(x)$  is a polynomial of degree  $n-1$ , then

$$\int_{-\infty}^{\infty} f(x) w(x) dx \approx \sum_{i=1}^n w_i f(x_i) \quad (1a)$$

weight function

Legendre  
 $w(x) = \begin{cases} 1 & \text{for } -1 < x < 1 \\ 0 & \text{elsewhere} \end{cases}$   
 Hermite  
 $w(x) = e^{-\frac{x^2}{2}}$

where  $w_1, \dots, w_n$  satisfies

$$\left. \begin{array}{l} f(x) = 1 \Rightarrow \int w(x) dx = \sum w_i \\ f(x) = x \Rightarrow \int x w(x) dx = \sum x_i w_i \\ \vdots \\ f(x) = x^{n-1} \Rightarrow \int x^{n-1} w(x) dx = \sum x_i^n w_i \end{array} \right\} n \text{ eqs.}$$

or

$$\left[ \begin{array}{ccc} 1 & 1 & 1 \\ x_1 & x_2 & x_3 \\ x_1^2 & x_2^2 & \dots \end{array} \right] \left[ \begin{array}{c} w_1 \\ w_2 \\ \vdots \\ w_n \end{array} \right] = \left[ \begin{array}{c} \int w(x) dx \\ \int x w(x) dx \\ \vdots \\ \int x^{n-1} w(x) dx \end{array} \right] \quad (1b)$$

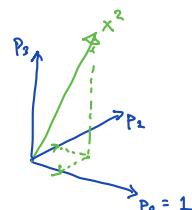
How to best choose evaluation points  $x_1, \dots, x_n$ ?

Choose  $\langle p, q \rangle = \int p(x) q(x) w(x) dx$  as inner product and construct orthogonal polynomials (basis vectors) from the polynomials  $1, x, x^2, \dots$  via Gram-Schmidt:

$$p_0(x) = 1$$

$$p_1(x) = x - \frac{\langle x, p_0 \rangle}{\langle p_0, p_0 \rangle} p_0(x) = \dots = x$$

$$p_2(x) = x^2 - \frac{\langle x^2, p_0 \rangle}{\langle p_0, p_0 \rangle} p_0(x) - \frac{\langle x^2, p_1 \rangle}{\langle p_1, p_1 \rangle} p_1(x) = \dots = x^2 - 1$$



$$P_3(x) = \dots = x^3 - 3x$$

for  $w(x) = e^{-\frac{x^2}{2}}$ ,  $P_0(x), P_1(x), \dots$  are known as Hermite polynomials.

Suppose that  $f(x)$  instead is of (at most) degree  $2n-1$ . Via polynomial long division by  $P_n(x)$  we have

$$f(x) = \underbrace{q(x) P_n(x)}_{\substack{\text{quotient polynomial} \\ \text{of order } \leq n-1}} + \underbrace{r(x)}_{\substack{\text{remainder} \\ \text{polynomial} \\ \text{of order } \leq n-1}} \quad (2)$$

E.g.  
 $x^3 - 2x^2 - 4 = \underbrace{(x^2 + x + 3)}_{\substack{\text{quotient } q(x)}} \underbrace{(x-3)}_{\substack{\text{divisor } r(x)}} + 5$

and hence

$$\int f(x) w(x) dx = \underbrace{\int q(x) P_n(x) w(x) dx}_{= \langle q, P_n \rangle = 0^*} + \int r(x) w(x) dx \quad (3)$$

since  $P_n$  is orthogonal to all polynomials spanned by  $P_{n-1}, P_{n-2}, \dots, P_0$  including  $q$  (of order  $n-1$ ).

Choosing the  $n$  real roots  $x_1, \dots, x_n$  of  $P_n(x)$  as evaluation (quadrature) points, (2) implies that  
 $r(x_i) = f(x_i)$  for  $i=1, \dots, n$ . (4)

Eqs. (1) ~ (4) then leads to

$$\int_{-\infty}^{\infty} f(x) w(x) dx \stackrel{(3)}{=} \int r(x) w(x) dx \stackrel{(2)}{=} \sum_{i=1}^n w_i r(x_i) \stackrel{(4)}{=} \sum_{i=1}^n w_i f(x_i)$$

Note: Exact as long as  $f(x)$  is a polynomial of at most order  $2n-1$ .

# Adaptive Gauss-Hermite quadrature (Lin & Pierce, 1994)

Marginal likelihood, clustered data

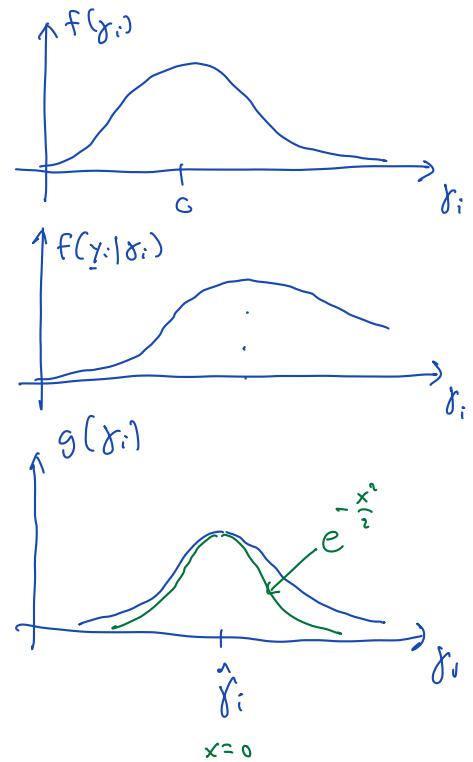
$$L(\beta, \theta) = f(\underline{y}_1, \underline{y}_2, \dots, \underline{y}_m)$$

$$= \prod_{i=1}^m f(\underline{y}_i)$$

$$= \prod_{i=1}^m \left( \int_{\mathbb{R}^n} \left( \prod_{j=1}^{n_i} f(y_{ij} | \gamma_i) \right) f(\gamma_i) d\gamma_i \right)$$

approx. Gaussian  
 if  $n_i$  is large

$\approx g(\gamma) \text{ approx. Gaussian}$



Letting  $\hat{\gamma} = \operatorname{argmax}_{\gamma} g(\gamma)$ ,  $h = -\frac{\partial^2}{\partial \gamma^2} \ln g(\gamma) \Big|_{\gamma=\hat{\gamma}}$  and  $x = \frac{\gamma - \hat{\gamma}}{1/\sqrt{h}}$ , such that  $d\gamma = \sqrt{h} dx$  and  $\gamma = \hat{\gamma} + \frac{x}{\sqrt{h}}$ , each integral

$$\begin{aligned}
 \int g(\gamma) d\gamma &= \sqrt{h} \int g\left(\hat{\gamma} + \frac{x}{\sqrt{h}}\right) dx \\
 &= \int f(x) e^{-\frac{x^2}{2}} dx \quad \text{where } f(x) = \frac{\sqrt{h} g\left(\hat{\gamma} + \frac{x}{\sqrt{h}}\right)}{e^{-x^2/2}} \left. \begin{array}{l} \text{deviation} \\ \text{of } g \text{ from} \\ \text{Gaussianity} \end{array} \right\} \\
 &\stackrel{nAQP}{\approx} \sum_{i=1}^n w_i f(x_i) \tag{5}
 \end{aligned}$$

Used by `lme4::glmer()`, `nAQP = ...`

Note that for  $nAQP = 1$ ,  $w_1 = \int e^{-\frac{x^2}{2}} dx = \sqrt{2\pi}$ ,  $x_1 = 0$  (1b) and (5) simplifies to  $\sqrt{2\pi h} g(\hat{\gamma})$ , the Laplace approximation of

$$\int g(\gamma) d\gamma$$



# Summary (November 16)

LMs

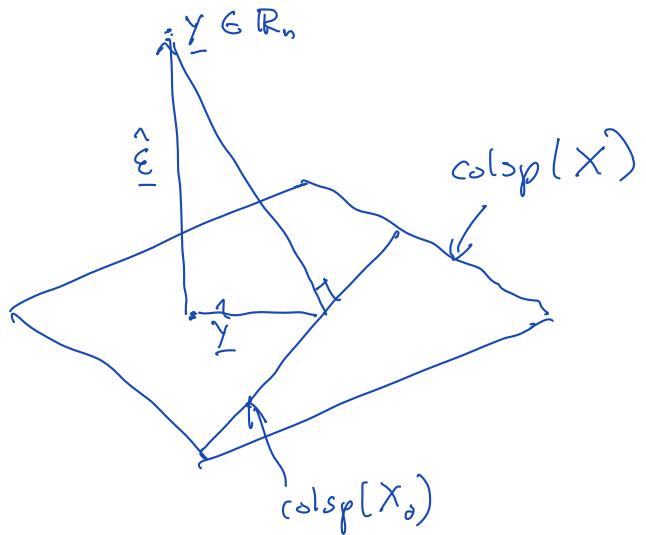
$$\underline{y} = \underline{X}\underline{\beta} + \underline{\varepsilon}, \quad \underline{\varepsilon} \sim N(0, \sigma^2 \mathbb{I}_n)$$

$$= x_1\beta_1 + x_2\beta_2 + \dots + x_p\beta_p + \underline{\varepsilon}$$

$$\hat{\underline{y}} = \arg \min_{\hat{\underline{y}} \in \text{colsp}(\underline{X})} \|\hat{\underline{y}} - \underline{y}\|$$

$$= \text{proj}_{\text{colsp}(\underline{X})} \underline{y}$$

$$= H\underline{y} = \underline{X} \underbrace{(\underline{X}^\top \underline{X})^{-1} \underline{X}^\top}_{\hat{\underline{\beta}}_{OLS}} \underline{y}$$



General linear model

$$\underline{y} = \underline{X}\underline{\beta} + \underline{\varepsilon}, \quad \underline{\varepsilon} \sim N(0, V)$$

↑

$$\underline{y}^* = \underline{X}^* \underline{\beta} + \underline{\varepsilon}^*, \quad \underline{\varepsilon}^* \sim N(0, I) \quad \underline{y}^* = V^{-1/2} \underline{y}$$

$$\hat{\underline{\beta}}_{GLS} = (\underline{X}^\top V^{-1} \underline{X})^{-1} \underline{X}^\top V \underline{y}$$

LMM

$$\underline{Y} = \underline{X}\beta + \underline{U}\gamma + \underline{\varepsilon}, \quad \gamma \sim N(0, G), \quad \varepsilon \sim N(0, R)$$

For example, random slope, intercept model

$$Y_{ij} = (\beta_0 + \gamma_{0,i}) + (\beta_1 + \gamma_{1,i}) X_{ij} + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim N(0, \sigma^2)$$

for  $i = 1, 2, \dots, m$ ,  $j = 1, 2, \dots, n_i$

$$\begin{bmatrix} \gamma_{0,i} \\ \gamma_{1,i} \end{bmatrix} \stackrel{iid}{\sim} N\left(0, \begin{bmatrix} T_0^2 & T_{01} \\ T_{01} & T_1^2 \end{bmatrix}\right)$$

ML:

$$\ell_p(\theta) = \max_{\beta} \ell(\beta, \theta) = \ell(\hat{\beta}(\theta), \theta)$$

$$\hat{\beta}(\theta) = \left[ \underline{X}^\top \underline{V}^{-1}(\theta) \underline{X} \right]^{-1} \underline{X}^\top \underline{V}^{-1}(\theta) \underline{y} \quad \text{where} \quad V(\theta) = U G(\theta) U^\top + R(\theta)$$

REML:

$$\begin{aligned} L_R(\theta) &\stackrel{\text{def}}{=} \int L(\beta, \theta) d\beta \\ &= \left| \underline{X}^\top \underline{X} \right|^{\frac{-1}{2}} f_{\tilde{w}} \left( \underbrace{\underline{A}^\top \underline{y}}_{(n-p) \times n} \right) \\ &= L_p(\theta) \left| \underline{X}^\top \underline{V}(\theta)^{-1} \underline{X} \right|^{\frac{1}{2}} \end{aligned}$$

GLMs. For  $i=1, 2, \dots, n$

$$y_i \sim \text{bin}(n_i, \pi_i)$$

Link functions  $\eta_i = g(\pi_i)$  Response function

$$\text{logit } \pi_i = \ln \frac{\pi_i}{1-\pi_i}$$

$$\pi_i = \frac{1}{1+e^{-\eta_i}}$$

$$\text{probit } \pi_i = \Phi^{-1}(\pi_i)$$

$$\pi_i = \Phi(\eta_i)$$

$$\text{cloglog} \approx \ln(-\ln(1-\pi_i))$$

$$\pi_i = 1 - e^{-e^{\eta_i}}$$

$$y_i \sim \text{Poisson}(\lambda_i)$$

$$\log(\lambda_i) = \eta_i$$

$$\lambda_i = e^{\eta_i}$$

cloglog examples

~ Presence-absence data

Poisson process:  $z_i \sim \text{Poisson}(\lambda_i)$ ,  $\ln \lambda_i = \eta_i$

Suppose we only observe

$$\pi_i = P(y_i = 1) = P(z_i \geq 1)$$

$$= 1 - P(z_i = 0)$$

$$= 1 - \frac{e^{-\lambda_i} \lambda_i^0}{0!} = 1 - e^{-\lambda_i}$$

$$= 1 - e^{-e^{y_i}}$$

- Survival data :  $T_i \sim \text{Weibull}$

$$\pi_i = P(\text{Dead after a time interval of length } t_i)$$

$$= P(T_i < t_i)$$

$$= 1 - e^{-e^{-at_i^b}}$$

$y_i = \beta_0 + \beta_1 \ln t_i$   
 $\ln a + b \ln t_i$

$$= 1 - e^{-e^{-\ln a - b \ln t_i}}$$

Probit link

$$\text{Age of monarche } T_i \sim N(\mu, \sigma^2)$$

$$\pi_i = P(y_i = 1) = P(T_i \leq \underline{t_i})$$

$$= \Phi\left(\frac{t_i - \mu}{\sigma}\right)$$

$$= \Phi\left(-\underbrace{\frac{\mu}{\sigma}}_{\beta_0} + \underbrace{\frac{1}{\sigma} t_i}_{\beta_1}\right)$$

# Multinomial models

$$\underline{Y_i} = \begin{bmatrix} Y_{i1} \\ Y_{i2} \\ \vdots \\ Y_{ic} \end{bmatrix} \sim \text{Multinomial} \left( n_i, \underline{\pi}_i \right) \quad \underline{\pi}_i = \begin{bmatrix} \pi_{i1} \\ \pi_{i2} \\ \vdots \\ \pi_{ic} \end{bmatrix}$$

Nominal models

$$\pi_{ir} = \frac{e^{\underline{x}_i^T \beta_r}}{1 + \sum_{s=1}^c e^{\underline{x}_i^T \beta_s}} = P(u_{ir} - \max_{s \neq r} u_{is} < 0)$$

↑  
non-parallel slopes

$$u_{ir} = \underline{x}_i^T \beta_r + \varepsilon_{ir}$$

↑  
~Gumbel

Ordinal models (example tomorrow)

$$\text{logit}(Y \leq r) = \theta_r + \underline{x}_i^T \underline{\beta}$$

↑  
parallel slopes

# GLMMs

Clustered data (unlike problem 2 in project 3)

$$y_{ij} | \underline{\gamma}_i \sim \text{bin}(n_{ij}, \pi_{ij})$$

or

$$y_{ij} | \underline{\gamma}_i \sim \text{Poisson}(\lambda_{ij})$$

and

$$\underline{\gamma}_i \sim N(\underline{\Omega}, Q)$$

Marginal likelihood

$$L(\beta, \theta) = f(\underline{y})$$

LTP

$$= \int f(\underline{y}, \underline{\gamma}) d\underline{\gamma}$$

$$= \int f(\underline{y} | \underline{\gamma}) f(\underline{\gamma}) d\underline{\gamma}$$

Non-Gaussian      In  
                        Gaussian

$\approx$  Laplace approximation

$\approx$  Gauss-Hermite quadrature

Difference between marginal and conditional model

### Possible forms of mathematical notation

Symbolic (Wilkinson & Rogers) notation

$$\sim 1 + x + j + x_i^i$$

or

$$\sim \underset{\substack{\uparrow \\ \text{numeric}}}{\text{age}} + \underset{\substack{\uparrow \\ \text{factor} \\ (\text{with } k \\ \text{levels})}}{\text{sex}} + \underset{\substack{\curvearrowleft \\ \text{interaction}}}{\text{age} \cdot \text{sex}}$$

Option 1: For  $i = 1, 2, \dots, n$

$$y_i = \mu + \beta \cdot x_i + \alpha_{j(i)} + \gamma_{j(i)} x_i \quad \begin{pmatrix} \text{see e.g. Wood} \\ \text{pp. 194 - 195} \end{pmatrix}$$

Option 1b:

$$y_i = \mu + \beta \cdot \text{age}_i + \alpha_{\text{sex}(i)} + \gamma_{\text{sex}(i)} \cdot \text{age}_i$$

Option 2: For  $j = 1, 2, \dots, k$  and  $i = 1, 2, \dots, n_j$

$$y_{ij} = \mu + \beta \cdot x_{ij} + \alpha_j + \gamma_j \cdot x_{ij}$$

Option 3:  $\downarrow \text{boldsymbol}\{\beta\}$

$$\underline{y} = X \underline{\beta}$$

If we have two levels  $k \geq 2$

$$= \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad ? \quad \begin{bmatrix} \mu \\ \beta \\ \alpha_2 \\ \gamma_2 \end{bmatrix}$$

### Project 3

$y_{ijk}$ : Number of goals of team  $i$  against team  $j$   
when playing away ( $k=2$ ) or at home ( $k=1$ )

Model: For  $i=1, \dots, 16$ ,  $j=1, \dots, 16$ ,  $i \neq j$  and  $k=1, 2$

$$y_{ijk} | \beta_i, \gamma_j \sim \text{Poisson}(\lambda_{ijk})$$

where

$$\ln \lambda_{ijk} = \underbrace{\mu + \alpha_k}_{\substack{\text{fixed effects} \\ \text{away vs. home}}} + \underbrace{\beta_i}_{\substack{\text{attacking team}}} + \underbrace{\gamma_j}_{\substack{\text{defending team}}}$$

where

$$\beta_i \stackrel{iid}{\sim} N(0, \tau_\beta^2) \quad \text{and} \quad \gamma_j \stackrel{iid}{\sim} N(0, \tau_\gamma^2)$$

Alternatively; For  $i=1, 2, \dots, 480$   
 $y_i \sim \text{Poisson}(\lambda_i)$

where

$$\ln \lambda_i = \mu + \alpha_{k(i)} + \beta_{j(i)} + \gamma_{l(i)}$$

Alternative 2:

$$\underline{y} \mid \underline{\gamma} \sim \text{Poisson}(\underline{\gamma})$$

$$\underline{\gamma} = X \beta + U \underline{\delta}, \quad \underline{\delta} \sim N(0, G)$$

TMA4315, 2006, problem 1

19 psychotic male ( $x_{i2}=1$ ) and female ( $x_{i2}=0$ ) patients with initial CGI =  $x_{i3}$  given injection every second week ( $x_{i1}=0$ ) or every third week ( $x_{i1}=1$ ).  $\underline{Y}_i = \begin{bmatrix} Y_{i1} \\ Y_{i2} \\ Y_{i3} \end{bmatrix}$  number of patients with final CGI = 0, 1, 2.

a) Formulate prop. odds model with no interaction

Response for  $i = 1, 2, \dots, 12$  is

$$\underline{Y}_i = (Y_{i1}, Y_{i2}, Y_{i3})^T \sim \text{Multinom}(\eta_i, \underline{\pi}_i)$$

where

$$\underline{\pi}_i = (\pi_{i1}, \pi_{i2}, \pi_{i3})^T$$

$C + 1 = 3$  categories

Covariates

$$\underline{x}_i = (x_{i1}, x_{i2}, x_{i3})^T$$

$\uparrow$              $\uparrow$              $\uparrow$   
 freq. of      sex      initial  
 medical                CGT  
 treatment

a) Proportional odds model

$$P(CGT_i \leq r) = P\left(-\underline{x}_i^T \beta + u_i \leq \theta_r\right)$$

$\sim$

$$\sim \text{logistic}(0, 1)$$

$$= P(u_i \leq \theta_r + \underline{x}_i^\top \beta)$$

$$= \text{logit}^{-1}(\theta_r + \underline{x}_i^\top \beta)$$

$$\approx \frac{1}{1 + e^{-(\theta_r + \underline{x}_i^\top \beta)}}$$

$$\theta_0 < \theta_1 < \theta_2 < \theta_3$$

$\downarrow$   
 $\infty$

Unknown parameters:  $\theta_1, \theta_2, \beta_1, \beta_2, \beta_3$

$$\pi_{ir} = P(CGI_i = r)$$

$$= P(CGI_i \leq r) - P(CGI_i \leq r-1)$$

$$\approx \frac{1}{1 + e^{-(\theta_r + \underline{x}_i^\top \beta)}} - \frac{1}{1 + e^{-(\theta_{r-1} + \underline{x}_i^\top \beta)}}$$

b) The model is that

$$P(CGI \leq r) = \text{logit}^{-1}(\theta_r + \underline{x}_i^\top \beta)$$

$$\text{log-odds } P(CGI \leq r) = \theta_r + \underline{x}_r^T \underline{\beta}$$

↑

$$\ln \frac{P(CGI \leq r)}{P(CGI > r)} = \theta_r + \underline{x}_r^T \underline{\beta}$$

$$\frac{P(CGI \leq r)}{P(CGI > r)} = e^{\theta_r + \underline{x}_r^T \underline{\beta}}$$

If  $x_{r3}$  increases by one unit the cumulative odds changes by a factor (an odds ratio) of  $e^{\beta_3}$ .

c) Saturated model

$$(\pi_{i1}, \pi_{i2}), i=1, 2, \dots, 12$$

free parameters ( $P_1 = 24$ )

with MLEs  $\hat{\pi}_{ij}^{(sat)} = \frac{y_{ij}}{n_i}$

The deviance is given by

$$D = 2 \left[ l\left(\hat{\pi}_{\text{sat}}\right) - l\left(\hat{\pi}_{\text{fitted}}\right) \right]$$

$$= 2 \sum_{i=1}^{12} \sum_{j=1}^3 Y_{ij} \ln \left( \frac{Y_{ij}/n_i}{\hat{\pi}_{ij}} \right)$$

The deviance (under  $H_0$ ) is

chi-square with

$$P_1 - P_0 = 24 - 5 = 19$$

according to asymptotic theory.

d) Model in symbolic notation:

$$\underbrace{B * K}_{+} + \underbrace{B * I}_{+}$$

$$\underbrace{B + K + B:K + I + B:I}_{+}$$

In mathematical notation the model is that

$$\text{logit}(\text{P}(CGI \leq r)) = \theta_r + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 \\ + \beta_4 x_1 x_2 + \beta_5 x_1 x_3$$

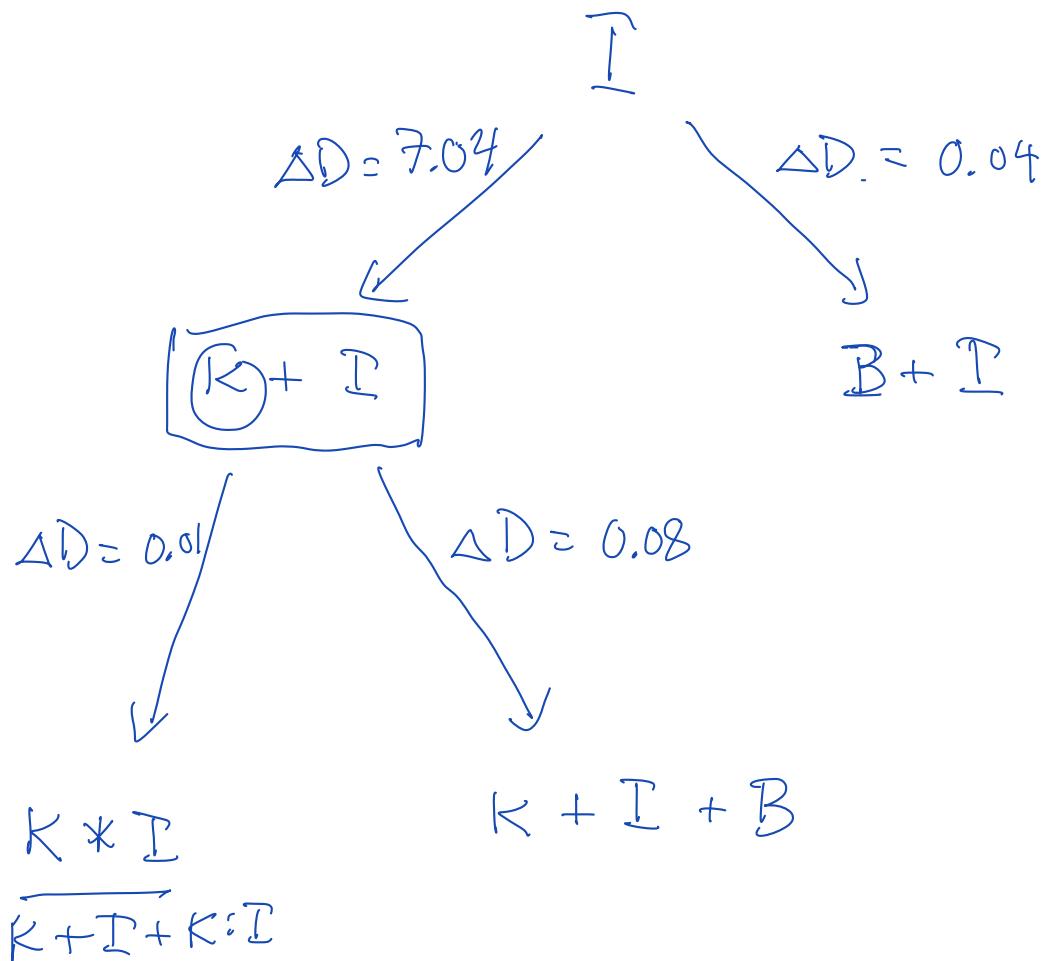
$\uparrow \quad \beta \quad \times \quad I$   
 $z$

This model has  $2 + 5 = 7$

The deviance for this model  $24 - 7 = 17$

Critical value

$$\chi^2_{0.05, 1} = 3.84$$



The best model (based on hypothesis testing) is  $K+I$

Could have done model selection via AIC.

e) Estimates of  $e^{\hat{\beta}_k}$ ,  $k=1, 2, 3$ .

$$\hat{e}^{\hat{\beta}_1} = e^{\hat{\beta}_1} = e^{-0.21} = 0.8026$$

Confidence interval for  $\hat{\beta}_1$  (based on

$$\frac{\hat{\beta}_1 - \beta_1}{\text{SE}(\hat{\beta}_1)} \stackrel{\text{approx}}{\sim} N(0, 1)$$

$$\hat{\beta}_1 \pm 1.96 \cdot \text{SE}(\hat{\beta}_1) = -0.21 \pm 1.96 \cdot 0.75$$
$$= (-1.81, 3.53)$$

Confidence interval for  $e^{\hat{\beta}_1}$

Recall that if  $(A, B)$  is a conf. int for  $\theta$ , then

$$P(A < \theta < B) = 1 - \alpha$$

and

$$P(g(A) < g(\theta) < g(B)) = 1 - \alpha$$

such that

$(g(A), g(B))$  is a  $(1-\alpha)$ -conf. int  
for  $g(\theta)$

$$\text{Hence, } (e^{0.18}, e^{3.53}) = (0.18, 3.53)$$

$$\hat{\pi}_{\hat{g}1} = P(CGI_{\hat{g}} = 1)$$

$$= P(CGI_{\hat{g}} \leq 1)$$

$$= \frac{1}{1 - e^{-\left(\hat{\theta}_1 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \hat{\beta}_3 x_3\right)}}$$

$\gamma_{ir}$

$$= \frac{1}{1 - e^{-\left(8.47 + -0.21 \cdot 0 + -2.15 \cdot 0 - 185 \cdot 5\right)}}$$

$r$

$\dots$

$$\text{Var} \left( \hat{\theta}_1 + \mathbf{x}_i^T \hat{\beta} \right)$$

$$= \text{Var} \left( \begin{bmatrix} 1 & 0 & \mathbf{x}_i^T \end{bmatrix} \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \vdots \\ \hat{\beta}_0 \end{bmatrix} \right)$$

$A$

$$= A \text{Var} \left( \begin{bmatrix} \hat{\theta} \\ \hat{\beta} \end{bmatrix} \right) A^T$$

||

$$F^{-1} \left( \begin{bmatrix} \hat{\theta} \\ \hat{\beta} \end{bmatrix} \right) \leftarrow \text{vcov}(\text{model})$$

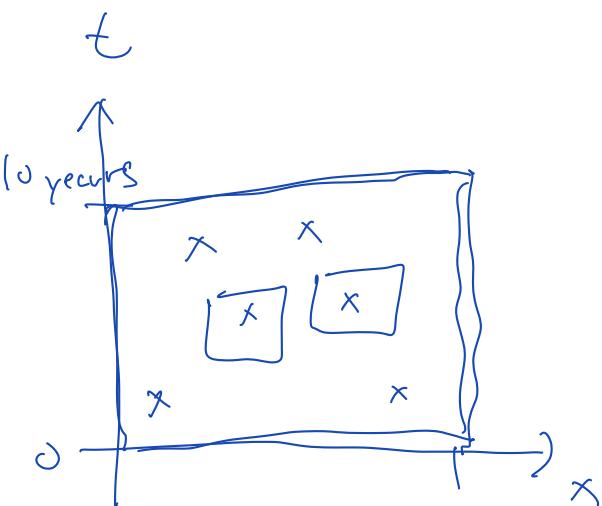
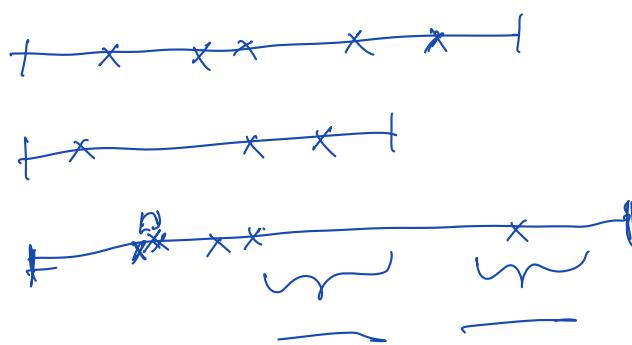
Delta method : Since  $\hat{\pi}_{ij1} = h(\hat{y}_{ij1})$

$$\text{Var}(\hat{\pi}_{ij1}) = \left( \frac{dh}{dy} \Big|_{y=\hat{y}_{ij1}} \right)^2 \text{Var}(\hat{\theta} + \hat{x}_j^\top \hat{\beta})$$

ST2304 2016, problem 3

a)

SFSSFFSss 



log-link : ensures that the poisson parameter  $> 0$

The model is

$$\ln \lambda_i = \underline{x}_i^T \underline{\beta} + \log(\text{length}_i)$$

$$\begin{aligned}\lambda_i &= e^{\underline{x}_i^T \underline{\beta}} \\ &= e^{\underline{x}_i^T \underline{\beta} + \log(\text{length}_i)}\end{aligned}$$

The offset term means that we assume direct proportionality between  $\lambda_i$  and  $\text{length}_i$ .

b) Changing the speed limit from  $80$  to  $70$  changes  $E\lambda_i = \lambda_i$  by a factor of

$$e^{\hat{\beta}_1 (\ln 70 - \ln 80)} = e^{\hat{\beta}_1 \ln \frac{70}{80}}$$

$$= \left(\frac{7}{8}\right)^{\hat{\beta}_1} = \left(\frac{7}{8}\right)^{4.92} = 0.51$$

i.e.  $49\%$  reduction in number of

accidents.

Absolute change greatest is for roads surrounded by woodland.

c) Overdispersion

- Positive covariance between subfactors

- Missing covariates

~~(Unexplained variation)~~

- Incorrect link function

Since  $D < \chi^2_{0.05, 295}$  we don't reject

$H_0$ : no overdispersion

$$qchisq(0.95, df=295) \approx 350$$

a) The model assumes that

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

and independent where

$$\ln \lambda_i = \beta_0 + \beta_1 x_i + \log(\text{area}_i)$$

where  $\beta_0$  and  $\beta_1$  are unknown parameters

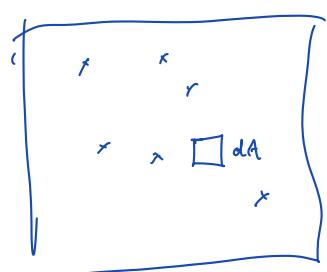
If altitude  $x_i$  increases by 1000m, the expected density ( $\text{plants} \cdot \text{m}^{-2}$ )

$$E\left(\frac{Y_i}{\text{area}_i}\right) = \frac{\lambda_i}{\text{area}_i} = e^{\beta_0 + \beta_1 x_i}$$

changes by a factor of  $e^{\beta_1 \cdot 1000} = e^{-0.58} = 0.56$ ,  
i.e., plant density is reduced by 44%.

b) Poisson assumption follows if

Individual plants occurrences  
in spatial Poisson process



Indep. between disjoint subareas  
Constant intensity  
 $P(\text{one occurrence in } dA) = \lambda dA$

Offset term:

$$\text{A priori } \lambda_i = E[Y_i] \propto \text{area}_i$$

$$\text{Model } \ln \lambda_i = \beta_0 + \beta_1 x_i + \ln(\text{area}_i)$$

①  $\beta_0 + \beta_1 x_i$

$$\lambda_i = \text{area}_i \cdot e$$

c) Overdispersion? Under  $H_0$  (no overdisp.)

$$D \sim \chi^2_{100-2}$$

↑  
98

Critical value

$$\chi^2_{0.95, 98} = 122.1$$

↑  
achisq

Obs. value

$$D = 235.$$

Reject  $H_0$  in favor of  $H_1$ : overdisp.

Mechanisms: Non-indep.

Missing covariates

Wrong link function

Estimated disp. par.  $\phi$  is

$$\hat{\phi} = \frac{D}{n-p} = \frac{235}{98} = 2.4$$

This changes Wald 2-test to Wald T-test  
with

$$T = \frac{\hat{\beta}_2}{\sqrt{\hat{\phi}} \text{SE}(\hat{\beta}_2)} = \frac{-0.00058}{\sqrt{2.4} \cdot 0.00045} = -0.909$$

Since  $|T| < t_{0.025, 98}$  we keep  $H_0: \beta_2 = 0$

d) Including year as a random effect, the model is  
that

$$Y_i | \gamma_{\text{year}(i)} \sim \text{Poisson}(\lambda_i)$$

where

$$\ln \lambda_i = \beta_0 + \beta_1 x_i + \gamma_{\text{year}(i)} + \ln(\text{area}_i)$$

and

$$\gamma_{\text{year}} \stackrel{\text{iid}}{\sim} N(0, \sigma_0^2) \quad \text{for year} = 2001, 2002, \dots, 2020.$$

Alternative notation (Fahrmeir)

$$Y_{ij} | \gamma_i \sim \text{Poisson}(\lambda_{ij})$$

where

$$\lambda_{ij} = \beta_0 + \beta_1 x_{ij} + \ln(\text{area}_{ij}) + \gamma_i$$

for years (clusters)  $i = 1, \dots, m$  and obs.  
 $j = 1, 2, \dots, n_i$

• MLE of  $\sigma_0^2 = \text{Var}(\gamma_{\text{year}})$  is  $\sigma_0^2 = 0.44$

• Given  $\gamma_{\text{year}} = 0$ , for  $x_i = 0$  (at sea level) and  $\text{area}_i = 50$

$$E(Y_i | \gamma_{\text{year}} = 0) = e^{\beta_0 + \ln(50)} = e^{-2.77} \cdot 50 = 3.13$$

Choosing a year at random

$$EY_i = E E(Y_i | \gamma)$$

$$= E e^{\beta_0 + \ln 50 + \gamma}$$

$$= 3.13 e^{\gamma} = 3.13 e^{0 + \sigma_0^2/2} = 3.90$$

e) Does the random year effect sign. improve model?

$$H_0: \sigma_0^2 = 0$$

( $H_0$  at boundary)

$$H_1: \sigma_0^2 > 0$$

- Under  $H_0$

$$E(\underline{\zeta}(\beta_1, \gamma_0)) = E\left(\left(\frac{\partial}{\partial \beta_0}, \frac{\partial}{\partial \beta_1}, \frac{\partial}{\partial \gamma_0}\right)^T \ell\right) = 0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

- $LRT = 2(\ell(H_1) - \ell(H_0)) \sim 0.5 \chi^2_0 : 0.5 \chi^2_1$

- Rejection region. We want

$$P(\text{Type I error}) = \alpha$$

$$P(LRT > c | H_0) = \alpha$$

$$\underbrace{\frac{1}{2} P(\chi^2_0 > c)}_{\approx 0} + \frac{1}{2} P(\chi^2_1 > c) = \alpha$$

$$P(\chi^2_1 > c) = 2\alpha$$

$$c = \chi^2_{2\alpha, 1} = \chi^2_{0.10, 1} = 6.63$$

↑  
tabell

Observed value

$$LRT = 2(-177.5 - (-225.1)) = 96$$

Concl.: Reject  $H_0$  in favour of  $H_1$

ST2304, 2013, problem 3

a) count ~ weekend + precip }  
 ||  
 yes, no (2 levels)

b) count ~ weekday + precip }  
 ||  
 mon, tue, ..., sun (7 levels)

$H_0$  is nested in  $H_1$  since  $H_0$  is equivalent

to  $H_1$  when  $\beta_{\text{mon}} = \beta_{\text{tue}} = \dots = \beta_{\text{fri}}$  and  $\beta_{\text{sat}} = \beta_{\text{sun}}$

Alternative view: Design matrices

$$X_0 = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ \dots \\ 1 & 0 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \quad X_1 = \begin{bmatrix} 1 & & & & \\ & 1 & 1 & & \\ & & 1 & 1 & \\ & & & 1 & 1 \\ & & & & 1 \end{bmatrix}$$

$$\Rightarrow \text{Colsp}(X_0) \subset \text{Colsp}(X_1)$$

$\Rightarrow H_0$  is nested in  $H_1$ .

Saturated model has  $X_s \in \mathbb{R}^n$

with  $\text{Colsp}(X_s) = \mathbb{R}^n$ .

All models nested in satur. mod.

ST2304, august 2015, problem 2 c, f, g

hemoglobin explained by

