

# Methods of Applied Statistics I

STA2101H F LEC9101

---

Week 8



October 29 2020

A promotional banner for an event. The background is a light brown color with a faint grid pattern. The text is white and black. Two circular portraits of speakers are on the right side. The overall design is clean and professional.

**Three Rs**  
— Reliability, Replicability, Reproducibility:  
the interplay between statistical science and data science

**Oct. 30 at 11 a.m.**

*The inaugural event of the Myles Hollander Distinguished Lectureship*

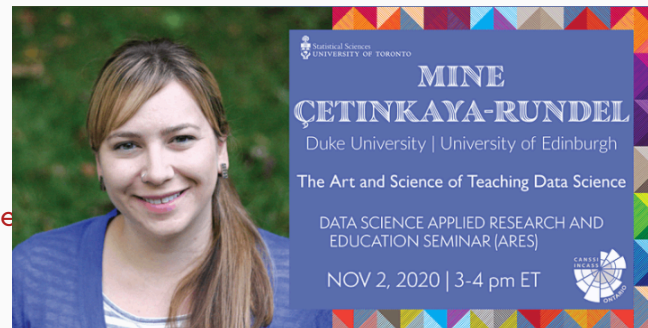
*Dr. Nancy Reid*

*Dr. Myles Hollander*

1. HW2 due November 5
  2. Measures of risk
  3. Modelling with binomial data (FELM §2.4–2.11)
  4. Generalized linear models (FELM Ch. 6)
  5. HW2 Questions
- November 2 3.00 – 4.00 Mine Çetinkaya-Rundel
  - [https://canssiontario.utoronto.ca/?mec-events=ares\\_cetinkaya-rundel\\_mine](https://canssiontario.utoronto.ca/?mec-events=ares_cetinkaya-rundel_mine)
  - “The art and science of teaching data science”

1. HW2 due November 5
2. Measures of risk
3. Modelling with binomial data (FELM §2.4–2.11)
4. Generalized linear models (FELM Ch. 6)
5. HW2 Questions

- November 2 3.00 – 4.00 Mine Çetinkaya-Rundel
- <https://canssiontario.utoronto.ca/?mec-evening>
- “The art and science of teaching data science”



# HALLOWEEN GGPLOT WORKSHOP!



Want to learn how to make compelling data visualizations with the powerful and flexible ggplot2 package in R?

Want an excuse to dress up for Halloween even if you're not leaving the house?

If your answer to one or both of those questions is "YES!" come join us for a very spooooooky workshop.

**Friday, Oct 30, 12:00–2:00 p.m. ET**

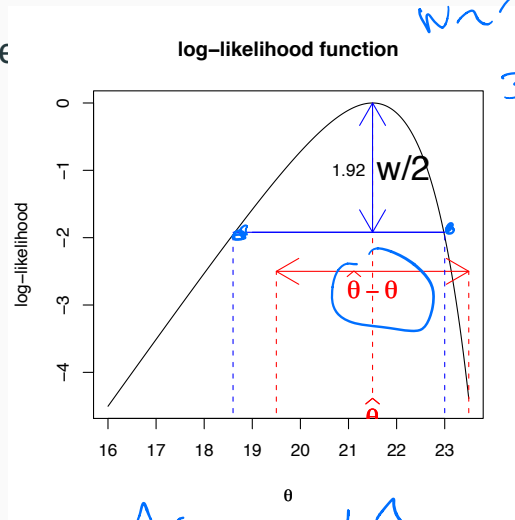
[Register here.](#)



- Regression for explanation
- observational data vs designed experiment; causality
- likelihood inference
- standardized maximum likelihood estimate (Wald test)
- likelihood ratio test
- modelling and inference for binary/binomial data
- saturated model and residual deviance
- interpretation of coefficients
- variable selection, residuals, diagnostics

*Wald test*

*$w \sim \chi^2_1$   
3.84*



$$\frac{\hat{\theta} - \theta}{\hat{\sigma}} \sim N(0, 1) = \text{Wald st.}$$

$$\hat{\theta} \pm 1.96 \cdot \hat{\sigma}$$

- see **posted handout** on case-control studies
- consider for simplicity binomial responses with a single binary covariate:

$$\log \frac{p_i}{1-p_i} = \text{logit}(p_i) \sim \beta_0 + \beta_1 z_i, \quad i = 1, \dots, n \quad y_i \sim \text{Bin}(n_i, p_i)$$

$y = 1$  "success"  
 $0$  "failure"

$$\hookrightarrow \begin{cases} \beta_0 & z_i = 0 \text{ "control"} \\ \beta_0 + \beta_1 & z_i = 1 \text{ "tut"} \end{cases}$$

$$\begin{aligned} \beta_1 \text{ effect of "tut" on } p_i &= \log\left(\frac{p_1}{1-p_1}\right) - \log\left(\frac{p_0}{1-p_0}\right) \\ &= \log\left(\frac{p_1}{1-p_1} \bigg/ \frac{p_0}{1-p_0}\right) \quad \left( e^{\hat{\beta}_1} = \frac{\hat{p}_1(1-\hat{p}_0)}{\hat{p}_0(1-\hat{p}_1)} \right) \text{ odds ratio} \end{aligned}$$

- no difference between groups  $\iff$  odds-ratio  $\equiv 1$

## ... Measures of risk

- we might be interested in **risk ratio**  $\frac{p_1}{p_0}$  instead of **odds ratio**  $\frac{p_1(1-p_0)}{p_0(1-p_1)}$
- also called **relative risk**

## ... Measures of risk

- we might be interested in **risk ratio**  $\frac{p_1}{p_0}$  instead of **odds ratio**  $\frac{p_1(1-p_0)}{p_0(1-p_1)}$
- also called **relative risk**
- if  $p_1$  and  $p_0$  are both small, ( $y = 1$  is rare), then

$$\frac{p_1}{p_0} \approx \frac{p_1(1-p_0)}{p_0(1-p_1)}$$

- sometimes  $p_1/p_0$  can be large but if  $p_1$  and  $p_0$  are both small the difference  $p_1 - p_0$  might also be very small



## ... Measures of risk

- we might be interested in **risk ratio**  $\frac{p_1}{p_0}$  instead of **odds ratio**  $\frac{p_1(1-p_0)}{p_0(1-p_1)}$
- also called **relative risk**
- if  $p_1$  and  $p_0$  are both small, ( $y = 1$  is rare), then

$$\frac{p_1}{p_0} \approx \frac{p_1(1-p_0)}{p_0(1-p_1)}$$

- sometimes  $p_1/p_0$  can be large but if  $p_1$  and  $p_0$  are both small the difference  $p_1 - p_0$  might also be very small
- in order to estimate the **risk difference** we need to know the baseline risk  $p_0$
- bacon sandwiches [www.youtube.com/watch?v=4szyEbU94ig](http://www.youtube.com/watch?v=4szyEbU94ig)
- risk calculator [realrisk.wintoncentre.uk/p8](http://realrisk.wintoncentre.uk/p8)

 1 / 1000

 3 / 1000 (2 extra cases)



Odds ratio 2.91; baseline risk 1/1000

# Biostats secret sauce

Whether we sample **prospectively** or **retrospectively**, the odds ratio is the same

	Lung cancer	
	1 cases	0 controls
smoke = 1 (yes)	688	650
smoke = 0 (no)	21	59
	709	709

$p_1 / (1 - p_1)$

"smoke / not"

$$\text{retro: OR} = \frac{(688/709)/(21/709)}{(650/709)/(59/709)} = \frac{688 \times 59}{650 \times 21} = 2.97$$

$$\text{prosp: OR} = \frac{\{688/(688 + 650)\} / \{21/(21 + 59)\}}{\{650/(688 + 650)\} / \{59/(21 + 59)\}} = \frac{688 \times 59}{650 \times 21} = 2.97$$

$p_1$  Lung c. vs not

see "case-control", FELM §2.5,6, SM §10.4.2

```
glm(cbind(r, m-r) ~ age + weight, data = mydata,
    family = binomial, link = logit )
```

```
?family
```

```
link
```

a specification for the model link function. This can be a name/expression, a literal character string, a length-one character vector, or an object of class "link-glm" ...

The gaussian family accepts the links (as names) identity, log and inverse; the binomial family the links logit, probit, cauchit, (corresponding to logistic, normal and Cauchy CDFs respectively) log and cloglog (complementary log-log)

def.

def.

$\ln(\eta_i^T \beta)$

hidden  
latent

$$\rightarrow z_i = x_i^T \gamma + \epsilon_i, \quad y_i = I(z_i \geq 0)$$

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

see also FELM Fig 2.3

$$y_i = 1 \text{ if } z_i \geq 0$$

$$p_i = P_n(z_i \geq 0) = 1 - \Phi\left(\frac{-x_i^T \gamma}{\sigma}\right) \quad \epsilon_i \geq -x_i^T \gamma$$

$$= \Phi\left\{ x_i^T \underbrace{\left(\frac{\gamma}{\sigma}\right)}_{\beta} \right\}$$

$$- \log(-\log(p_i)) = x_i^T \beta$$

$$\sum_i x_i^T \beta y_i + \dots$$

$$L(\beta) = \prod_{i=1}^n \binom{n_i}{y_i} \Phi(x_i^T \beta)^{y_i} \{1 - \Phi(\quad)\}^{n_i - y_i}$$

$$L(\beta) = \sum \log \{ \quad \}$$

ESTIMATION PROBLEMS

39

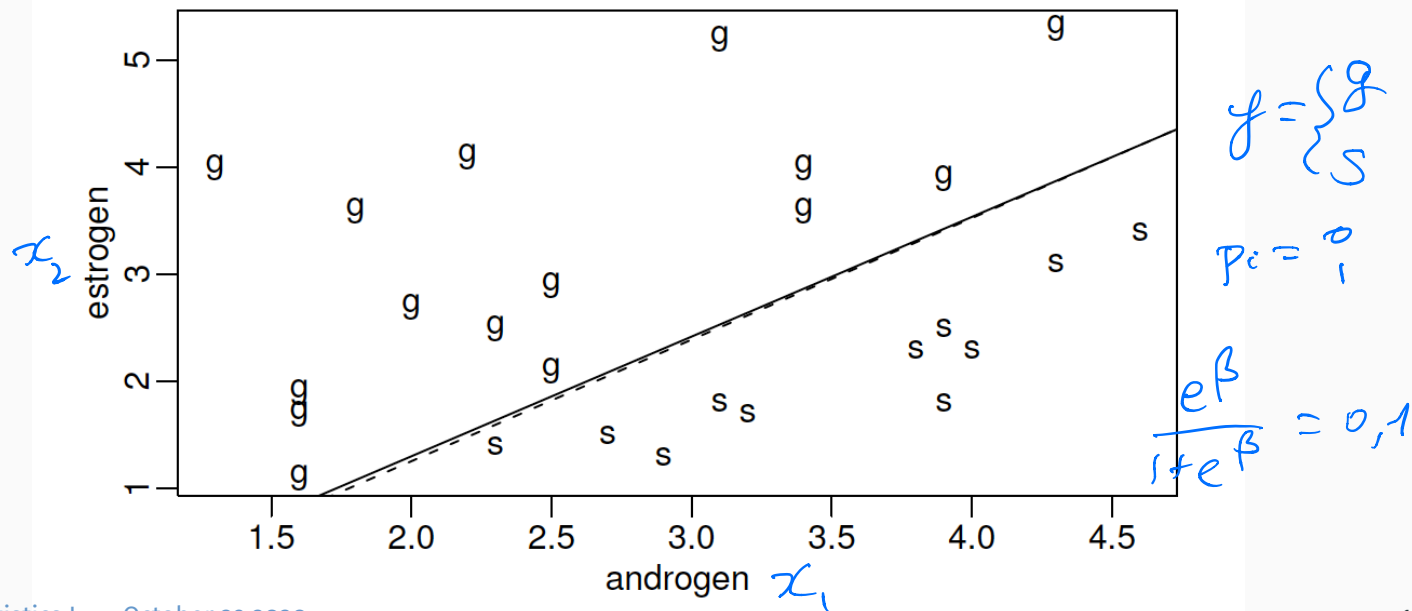


Figure 2.4 Levels of androgen and estrogen for 15 homosexual (g) and 11 heterosexual (s)

$$x_{1/2} = -\frac{\beta_0}{\beta_1}$$

$$\hat{x}_{1/2} = -\frac{\hat{\beta}_0}{\hat{\beta}_1}$$

$$\boxed{x^* + \beta}$$

$$\beta = (\beta_0, \beta_1)$$

n.v.

$$\text{var } g(\hat{\beta}) \doteq \underbrace{g'(\hat{\beta})^T}_{1 \times 2} \underbrace{\text{cov}(\hat{\beta})}_{2 \times 2} \underbrace{g'(\hat{\beta})}_{2 \times 1}$$

$$g'(\beta) = \begin{pmatrix} \frac{\partial}{\partial \beta_0} g \\ \frac{\partial}{\partial \beta_1} g \end{pmatrix}$$

$$g(\beta) = \frac{\beta_0}{\beta_1}$$

$$g'(\beta) = \begin{pmatrix} -\frac{1}{\beta_1} \\ \frac{\beta_0}{\beta_1^2} \end{pmatrix}$$

vcov

$$\hat{\text{var}}(\hat{x}_{1/2}) =$$

$$g(\beta) \approx g(\beta_0)$$

$$+ (\beta - \beta_0) g'(\beta_0)$$

$$\hat{p}(x^*) = \frac{e^{x^{*T}\hat{\beta}}}{1 + e^{x^{*T}\hat{\beta}}}$$

$$\text{var}(x^{*T}\hat{\beta}) = x^{*T} \text{cov}(\hat{\beta}) x^*$$

ilogit

$$x^{*T}\hat{\beta} \sim N(x^{*T}\beta, \underbrace{x^{*T} \text{cov}(\hat{\beta}) x^*}_{\text{vcov}})$$

95% CI for  $x^{*T}\beta$ :  $x^{*T}\hat{\beta} \pm 1.96 \sqrt{\text{vcov}} = (L, U)$

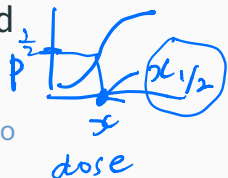
95% CI for  $\hat{p}(x^*) = (\text{ilogit}(L), \text{ilogit}(U)) \in (0, 1)$   $\hat{\beta} = \hat{\beta}(y)$   
m.l.e

$$\beta_0 + \beta_1 x = \text{logit}\{p(x)\}$$

$$y_i | x_i \sim \text{Bin}\{n_i, p(x_i)\}$$

ED50 and delta method

$$\beta_1 > 0$$



$$p = \frac{1}{2} \quad \log\left(\frac{1/2}{1-1/2}\right) = 0$$

$$\beta_0 + \beta_1 x_{1/2} = 0$$

$$x_{1/2} = -\frac{\beta_0}{\beta_1}$$



- $Y_i \sim \text{Bin}(n_i, p_i) \Rightarrow E(Y_i) = n_i p_i, \quad \text{Var}(Y_i) = n_i p_i (1 - p_i)$

residual deviance  $2 \{ \ell(\tilde{p}_i) - \ell(p_i(\hat{\beta})) \}$  sat. 'd

$= \dots \sum_{i=1}^n \{ \dots \}$  reg. model

$P = (\beta_1, \dots, \beta_p)$

$\text{var}(Y_i) = \phi n_i p_i (1 - p_i)$

$\phi$  "fudge factor"

overdispersion par.

$\tilde{\phi} = \frac{\text{Res. dev.}}{n-p}$  estimated

quasi-likelihood  $\begin{pmatrix} -\hat{\beta}_2 \\ \hat{s}_e^2 \hat{\beta}_q \end{pmatrix}$   $\sim N(\dots)$

overdisp.Rmd; overdisp.html

# Generalized linear models: theory

•  
p.m.f. for  
Bin, Po  
pdf for  
N, G, iG

$$f(y_i; \mu_i, \phi_i) = \exp\left\{ \frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i; \phi_i) \right\}$$

glm

- binom -
- norm
- Poiss -
- Gamma
- (- inv. G.)
- quasi-bin
- quasi-P.

# Generalized linear models: theory

- $f(y_i; \mu_i, \phi_i) = \exp\left\{\frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i; \phi_i)\right\}$
  - $E(y_i) = b'(\theta_i) = \mu_i$  defines  $\mu_i$  as a function of  $\theta_i$
- $\mu_i = E(y_i)$
- $\int y_i f(y_i) dy_i$

# Generalized linear models: theory

- $f(y_i; \mu_i, \phi_i) = \exp\left\{\frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i; \phi_i)\right\}$
- $E(y_i | x_i) = b'(\theta_i) = \mu_i$  defines  $\mu_i$  as a function of  $\theta_i$
- $g(\mu_i) = x_i^T \beta = \eta_i$  links the  $n$  observations together via covariates

$$\mu_i = E(y_i | x_i)$$
$$p_i = E(y_i) \quad \log \frac{p_i}{1-p_i}$$

# Generalized linear models: theory

•

$$f(y_i; \mu_i, \phi_i) = \exp\left\{\frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i; \phi_i)\right\}$$

•  $E(y_i | x_i) = b'(\theta_i) = \mu_i$  defines  $\mu_i$  as a function of  $\theta_i$

•  $g(\mu_i) = x_i^T \beta = \eta_i$  links the  $n$  observations together via covariates

•  $g(\cdot)$  is the **link** function;  $\eta_i$  is the **linear predictor**

$$\begin{array}{l} \underline{x}_i \quad p \times 1 \\ \underline{\beta} \quad p \times 1 \\ \hline \eta_i \quad p < n \end{array}$$

$$\begin{array}{l} \theta_i \\ \mu_i \\ \eta_i \end{array} \quad i=1, \dots, n$$

$$\underline{x}_i^T$$

$$\beta_j \quad j=1, \dots, p$$

# Generalized linear models: theory

- $$f(y_i; \mu_i, \phi_i) = \exp\left\{\frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i; \phi_i)\right\}$$
def = glm
- $E(y_i | x_i) = \underline{b'(\theta_i)} = \mu_i$  defines  $\mu_i$  as a function of  $\theta_i$ 
 $E(y_i) = \int y_i f(y_i) dy_i$
- $g(\mu_i) = x_i^T \beta = \eta_i$  links the  $n$  observations together via covariates
- $g(\cdot)$  is the **link** function;  $\eta_i$  is the **linear predictor**

- $\underline{\text{Var}(y_i | x_i)} = \phi_i \underline{b''(\theta_i)} = \phi_i \underline{V(\mu_i)}$ 
variance
prop's that follow

extendible  
Applied Statistics

$$g\{E(y_i)\} = x_i^T \beta$$

$$g(\mu_i) = x_i^T \beta$$

$$\text{var}(y_i | x_i) = \phi_i V(\mu_i)$$

$y_i$  ind't

# Generalized linear models: theory

•

$$f(y_i; \mu_i, \phi_i) = \exp\left\{\frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i; \phi_i)\right\}$$

- $E(y_i | x_i) = b'(\theta_i) = \mu_i$  defines  $\mu_i$  as a function of  $\theta_i$
- $g(\mu_i) = x_i^T \beta = \eta_i$  links the  $n$  observations together via covariates
- $g(\cdot)$  is the **link** function;  $\eta_i$  is the **linear predictor**
- $\text{Var}(y_i | x_i) = \phi_i b''(\theta_i) = \phi_i V(\mu_i)$
- $V(\cdot)$  is the **variance function**  $\text{var} = \sigma^2$ .

$$b(\mu_i) = \frac{1}{2} \mu_i^2 \\ b' = \mu_i \quad b'' = 1$$



- Normal
- Binomial
- Poisson
- Gamma/E
- Inverse G

```
family {stats}
```

R Documentation

## Family Objects for Models

### Description

Family objects provide a convenient way to specify the details of the models used by functions such as `glm`. See the documentation for `glm` for the details on how such model fitting takes place.

### Usage

```
family(object, ...)  
  
binomial(link = "logit")  
gaussian(link = "identity")  
Gamma(link = "inverse")  
inverse.gaussian(link = "1/mu^2")  
poisson(link = "log")  
quasi(link = "identity", variance = "constant")  
quasibinomial(link = "logit")  
quasipoisson(link = "log")
```

# Examples

• Normal:  $f(y_i; \mu_i, \sigma^2) = \frac{1}{\sqrt{(2\pi)\sigma}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \mu_i)^2\right\}$

$= \exp\left\{\frac{y_i\mu_i - (1/2)\mu_i^2}{\sigma^2} - (1/2)\log\sigma^2 - y_i^2/2\sigma^2 - (1/2)\log\sqrt{(2\pi)}\right\}$

$e^{\frac{y_i\theta_i - b(\theta_i)}{\phi_i} + c(y_i; \phi_i)}$

$\phi_i = \sigma^2, \quad \theta_i = \mu_i, \quad b(\mu_i) = \mu_i^2/2\sigma^2$

note  $b''(\mu_i) = 1$

$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(y_i - \mu_i)^2} = e^{-\frac{1}{2\sigma^2}y_i^2 + \frac{1}{\sigma^2}y_i\mu_i - \frac{1}{2\sigma^2}\mu_i^2 - \frac{1}{2}\ln\sigma^2 - \frac{1}{2}k_2\pi}$

$= e^{\frac{(y_i\mu_i - \mu_i^2/2)}{\sigma^2} + c}$

$E(y_i) = \mu_i = x_i^T \beta$

$\text{var}(y_i) = \sigma^2 = \phi_i$

$\phi_i \equiv \sigma^2$

$\theta_i \equiv \mu_i$  identity

$b'(\theta_i) = b'(\mu_i) = 2\mu_i/2 = \mu_i = E(y_i)$

# Examples

• Normal:  $f(y_i; \mu_i, \sigma^2) = \frac{1}{\sqrt{(2\pi)\sigma}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \mu_i^2)\right\}$

$(m_i \sigma^2)$   
link.  
=  $\exp\left\{\frac{y_i \mu_i - (1/2)\mu_i^2}{\sigma^2} - (1/2) \log \sigma^2 - y_i^2/2\sigma^2 - (1/2) \log \sqrt{(2\pi)}\right\}$

overdisp.  
 $\phi_i = \phi/m_i$   
quasi  
note  $b''(\mu_i) = 1$

$\phi_i = \sigma^2, \theta_i = \mu_i, b(\mu_i) = \mu_i^2/2\sigma^2$

• Binomial:  $f(r_i; p_i) = \binom{m_i}{r_i} p_i^{r_i} (1-p_i)^{m_i-r_i}; y_i = r_i/m_i$

$\mu_i = E\left(\frac{r_i}{m_i}\right) = \exp[m_i y_i \log\{p_i/(1-p_i)\} + m_i \log(1-p_i) + \log\left(\binom{m_i}{m_i y_i}\right)]$

known  
 $\phi_i = 1/m_i; \theta_i = \log\{p_i/(1-p_i)\}, b(p_i) = -\log(1-p_i)$

$\frac{1}{\phi_i} \theta_i = y_i$   
 $\frac{1}{\phi_i} b(\theta_i)$

$c(y_i, \phi_i)$   
 $\phi_i b''(\theta_i) = \text{var}(y_i)$   
Note  $p_i = \mu_i = E(y_i)$

$f(y_i) = \binom{m_i}{m_i y_i} p_i^{m_i y_i} (1-p_i)^{m_i - m_i y_i}$

$= \exp\left[m_i \left(\log \frac{p_i}{1-p_i}\right) y_i + m_i \log(1-p_i)\right]$

# Examples

- Normal:  $f(y_i; \mu_i, \sigma^2) = \frac{1}{\sqrt{(2\pi)\sigma}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \mu_i^2)\right\}$   
 $= \exp\left\{\frac{y_i\mu_i - (1/2)\mu_i^2}{\sigma^2} - (1/2)\log\sigma^2 - y_i^2/2\sigma^2 - (1/2)\log\sqrt{(2\pi)}\right\}$

$$\phi_i = \sigma^2, \quad \theta_i = \mu_i, \quad b(\mu_i) = \mu_i^2/2\sigma^2$$

note  $b''(\mu_i) = 1$

- Binomial:  $f(r_i; p_i) = \binom{m_i}{r_i} p_i^{r_i} (1 - p_i)^{m_i - r_i}; \quad y_i = r_i/m_i$   
 $= \exp[m_i y_i \log\{p_i/(1 - p_i)\} + m_i \log(1 - p_i) + \log\left(\binom{m_i}{m_i y_i}\right)]$

$$\phi_i = 1/m_i, \quad \theta_i = \log\{p_i/(1 - p_i)\}, \quad b(p_i) = -\log(1 - p_i)$$

Note  $p_i = \mu_i = E(y_i)$

- ELM (p.115) uses  $a_i(\phi)$  in place of  $\phi_i$ , later (p.117)  $a_i(\phi) = \phi/w_i$ ; later (p.118)  $w_i$  used for weights in IRWLS algorithm; SM uses  $\phi_i$ , later (p. 483)  $\phi_i = \phi a_i$

- $\ell(\beta; \mathbf{y}) = \sum \left\{ \frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i, \phi_i) \right\}$

- $\ell(\beta; \mathbf{y}) = \sum \left\{ \frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i, \phi_i) \right\}$
- $\mathbf{b}'(\theta_i) = \mu_i; \quad \mathbf{g}(\mu_i) = \mathbf{g}(\mathbf{b}'(\theta_i)) = \eta_i = \mathbf{x}_i^T \beta$

- $\ell(\beta; \mathbf{y}) = \sum \left\{ \frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i, \phi_i) \right\}$
- $\mathbf{b}'(\theta_i) = \mu_i; \quad \mathbf{g}(\mu_i) = \mathbf{g}(\mathbf{b}'(\theta_i)) = \eta_i = \mathbf{x}_i^T \beta$
- $\frac{\partial \ell(\beta; \mathbf{y})}{\partial \beta_j} = \sum \frac{\partial \ell_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \beta_j} = \sum \frac{y_i - \mathbf{b}'(\theta_i)}{\phi_i} \frac{\partial \theta_i}{\partial \beta_j}$

- $\ell(\beta; \mathbf{y}) = \sum \left\{ \frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i, \phi_i) \right\}$
- $\mathbf{b}'(\theta_i) = \mu_i; \quad \mathbf{g}(\mu_i) = \mathbf{g}(\mathbf{b}'(\theta_i)) = \eta_i = \mathbf{x}_i^T \beta$
- $\frac{\partial \ell(\beta; \mathbf{y})}{\partial \beta_j} = \sum \frac{\partial \ell_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \beta_j} = \sum \frac{y_i - \mathbf{b}'(\theta_i)}{\phi_i} \frac{\partial \theta_i}{\partial \beta_j}$
- $\mathbf{g}'(\mathbf{b}(\theta_i)) \mathbf{b}''(\theta_i) \frac{\partial \theta_i}{\partial \beta_j} = \mathbf{x}_{ij} = \mathbf{g}'(\mu_i) \mathbf{V}(\mu_i)$

See Slide 2



- $\ell(\beta; \mathbf{y}) = \sum \left\{ \frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i, \phi_i) \right\}$
- $\mathbf{b}'(\theta_i) = \mu_i; \quad \mathbf{g}(\mu_i) = \mathbf{g}(\mathbf{b}'(\theta_i)) = \eta_i = \mathbf{x}_i^T \beta$
- $\frac{\partial \ell(\beta; \mathbf{y})}{\partial \beta_j} = \sum \frac{\partial \ell_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \beta_j} = \sum \frac{y_i - \mathbf{b}'(\theta_i)}{\phi_i} \frac{\partial \theta_i}{\partial \beta_j}$
- $\mathbf{g}'(\mathbf{b}(\theta_i)) \mathbf{b}''(\theta_i) \frac{\partial \theta_i}{\partial \beta_j} = \mathbf{x}_{ij} = \mathbf{g}'(\mu_i) \mathbf{V}(\mu_i)$
- $\frac{\partial \ell(\beta; \mathbf{y})}{\partial \beta_j} = \sum \frac{y_i - \mu_i}{\phi_i \mathbf{g}'(\mu_i) \mathbf{V}(\mu_i)} \mathbf{x}_{ij} = \sum \frac{y_i - \mu_i}{a_i \mathbf{g}'(\mu_i) \mathbf{V}(\mu_i)} \mathbf{x}_{ij}$

See Slide 2

when  $\phi_i = a_i \phi$

- $\ell(\beta; \mathbf{y}) = \sum \left\{ \frac{y_i \theta_i - b(\theta_i)}{\phi_i} + c(y_i, \phi_i) \right\}$
- $b'(\theta_i) = \mu_i; \quad g(\mu_i) = g(b'(\theta_i)) = \eta_i = \mathbf{x}_i^T \beta$
- $\frac{\partial \ell(\beta; \mathbf{y})}{\partial \beta_j} = \sum \frac{\partial \ell_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \beta_j} = \sum \frac{y_i - b'(\theta_i)}{\phi_i} \frac{\partial \theta_i}{\partial \beta_j}$
- $g'(b(\theta_i)) b''(\theta_i) \frac{\partial \theta_i}{\partial \beta_j} = \mathbf{x}_{ij} = g'(\mu_i) \mathbf{V}(\mu_i)$
- $\frac{\partial \ell(\beta; \mathbf{y})}{\partial \beta_j} = \sum \frac{y_i - \mu_i}{\phi_i g'(\mu_i) \mathbf{V}(\mu_i)} \mathbf{x}_{ij} = \sum \frac{y_i - \mu_i}{a_i g'(\mu_i) \mathbf{V}(\mu_i)} \mathbf{x}_{ij}$

See Slide 2

when  $\phi_i = a_i \phi$

- matrix notation:

$$\frac{\partial \ell(\beta)}{\partial \beta} = \mathbf{X}^T \mathbf{u}(\beta), \quad \mathbf{X} = \frac{\partial \eta}{\partial \beta^T}, \quad \mathbf{u} = (u_1, \dots, u_n)$$

## Scale parameter $\phi_i$

- in most cases, either  $\phi_i$  is known, or  $\phi_i = \phi a_i$ , where  $a_i$  is known

## Scale parameter $\phi_i$

- in most cases, either  $\phi_i$  is known, or  $\phi_i = \phi a_i$ , where  $a_i$  is known
- Normal distribution,  $\phi = \sigma^2$

## Scale parameter $\phi_i$

- in most cases, either  $\phi_i$  is known, or  $\phi_i = \phi a_i$ , where  $a_i$  is known
- Normal distribution,  $\phi = \sigma^2$
- Binomial distribution  $\phi_i = m_i^{-1}$

## Scale parameter $\phi_i$

- in most cases, either  $\phi_i$  is known, or  $\phi_i = \phi a_i$ , where  $a_i$  is known
- Normal distribution,  $\phi = \sigma^2$
- Binomial distribution  $\phi_i = m_i^{-1}$
- Gamma distribution,  $\phi = 1/\nu$

## Scale parameter $\phi_i$

- in most cases, either  $\phi_i$  is known, or  $\phi_i = \phi a_i$ , where  $a_i$  is known

- Normal distribution,  $\phi = \sigma^2$

- Binomial distribution  $\phi_i = m_i^{-1}$

- Gamma distribution,  $\phi = 1/\nu$

- $$\frac{\partial \ell(\beta; \mathbf{y})}{\partial \beta_j} = \sum \frac{y_i - \mu_i}{\phi_i g'(\mu_i) V(\mu_i)} x_{ij} = \sum \frac{y_i - \mu_i}{a_i g'(\mu_i) V(\mu_i)} x_{ij}$$

when  $\phi_i = a_i \phi$

## Scale parameter $\phi_i$

- in most cases, either  $\phi_i$  is known, or  $\phi_i = \phi a_i$ , where  $a_i$  is known

- Normal distribution,  $\phi = \sigma^2$

- Binomial distribution  $\phi_i = m_i^{-1}$

- Gamma distribution,  $\phi = 1/\nu$

- $$\frac{\partial \ell(\beta; \mathbf{y})}{\partial \beta_j} = \sum \frac{y_i - \mu_i}{\phi_i g'(\mu_i) V(\mu_i)} x_{ij} = \sum \frac{y_i - \mu_i}{a_i g'(\mu_i) V(\mu_i)} x_{ij}$$

when  $\phi_i = a_i \phi$

- if  $\theta_i = g(\mu_i)$  **canonical link**, then  $g'(\mu_i) = 1/V(\mu_i)$ , and

$$\sum \frac{y_i x_{ij}}{a_i} = \sum \frac{y_i \hat{\mu}_i x_{ij}}{a_i}$$



## Solving maximum likelihood equation

- Newton-Raphson:  $\ell'(\hat{\beta}) = \mathbf{0} \approx \ell'(\beta) + (\hat{\beta} - \beta)\ell''(\beta)$

defines iterative scheme

## Solving maximum likelihood equation

- Newton-Raphson:  $l'(\hat{\beta}) = \mathbf{0} \approx l'(\beta) + (\hat{\beta} - \beta)l''(\beta)$

defines iterative scheme

- $\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} - \{l''(\hat{\beta}^{(t)})\}^{-1}l'(\hat{\beta}^{(t)})$

# Solving maximum likelihood equation

- Newton-Raphson:  $\ell'(\hat{\beta}) = \mathbf{0} \approx \ell'(\beta) + (\hat{\beta} - \beta)\ell''(\beta)$

defines iterative scheme

- $\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} - \{\ell''(\hat{\beta}^{(t)})\}^{-1}\ell'(\hat{\beta}^{(t)})$

- Fisher scoring:  $-\ell''(\beta) \leftarrow \mathbf{E}\{-\ell''(\beta)\} = \mathbf{i}(\beta)$

many books use  $I(\beta)$

# Solving maximum likelihood equation

- Newton-Raphson:  $\ell'(\hat{\beta}) = \mathbf{0} \approx \ell'(\beta) + (\hat{\beta} - \beta)\ell''(\beta)$

defines iterative scheme

- $\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} - \{\ell''(\hat{\beta}^{(t)})\}^{-1}\ell'(\hat{\beta}^{(t)})$

- Fisher scoring:  $-\ell''(\beta) \leftarrow \mathbf{E}\{-\ell''(\beta)\} = \mathbf{i}(\beta)$

many books use  $I(\beta)$

- $\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} + \{\mathbf{i}(\hat{\beta}^{(t)})\}^{-1}\ell'(\hat{\beta}^{(t)})$

# Solving maximum likelihood equation

- Newton-Raphson:  $\ell'(\hat{\beta}) = \mathbf{0} \approx \ell'(\beta) + (\hat{\beta} - \beta)\ell''(\beta)$

defines iterative scheme

- $\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} - \{\ell''(\hat{\beta}^{(t)})\}^{-1}\ell'(\hat{\beta}^{(t)})$

- Fisher scoring:  $-\ell''(\beta) \leftarrow \mathbf{E}\{-\ell''(\beta)\} = \mathbf{i}(\beta)$

many books use  $\mathbf{I}(\beta)$

- $\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} + \{\mathbf{i}(\hat{\beta}^{(t)})\}^{-1}\ell'(\hat{\beta}^{(t)})$

- applied to matrix version:  $\mathbf{X}^T \mathbf{u}(\hat{\beta}) = \mathbf{0} \doteq \mathbf{X}^T \mathbf{u}(\beta) + (\hat{\beta} - \beta)\mathbf{X}^T \frac{\partial \mathbf{u}(\beta)}{\partial \beta^T}$

slide 5

# Solving maximum likelihood equation

- Newton-Raphson:  $\ell'(\hat{\beta}) = \mathbf{0} \approx \ell'(\beta) + (\hat{\beta} - \beta)\ell''(\beta)$

defines iterative scheme

- $\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} - \{\ell''(\hat{\beta}^{(t)})\}^{-1}\ell'(\hat{\beta}^{(t)})$

- Fisher scoring:  $-\ell''(\beta) \leftarrow \mathbf{E}\{-\ell''(\beta)\} = \mathbf{i}(\beta)$

many books use  $I(\beta)$

- $\hat{\beta}^{(t+1)} = \hat{\beta}^{(t)} + \{\mathbf{i}(\hat{\beta}^{(t)})\}^{-1}\ell'(\hat{\beta}^{(t)})$

- applied to matrix version:  $\mathbf{X}^T \mathbf{u}(\hat{\beta}) = \mathbf{0} \doteq \mathbf{X}^T \mathbf{u}(\beta) + (\hat{\beta} - \beta)\mathbf{X}^T \frac{\partial \mathbf{u}(\beta)}{\partial \beta^T}$

slide 5

- change to Fisher scoring:  $\hat{\beta} = \beta + \mathbf{i}(\beta)^{-1}\mathbf{X}^T \mathbf{u}(\beta)$

## ... maximum likelihood equation

- $\hat{\beta} = \beta + \mathbf{i}(\beta)^{-1} \mathbf{X}^T \mathbf{u}(\beta)$

$$\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} = \sum \frac{-b''(\theta_i)}{\phi_i} \left( \frac{\partial \theta_i}{\partial \beta_j} \right) \left( \frac{\partial \theta_i}{\partial \beta_k} \right) + \sum \frac{y_i - b'(\theta_i)}{\phi_i} \frac{\partial^2 \theta_i}{\partial \beta_j \partial \beta_k}$$

## ... maximum likelihood equation

- $\hat{\beta} = \beta + \mathbf{i}(\beta)^{-1} \mathbf{X}^T \mathbf{u}(\beta)$

$$\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} = \sum \frac{-b''(\theta_i)}{\phi_i} \left( \frac{\partial \theta_i}{\partial \beta_j} \right) \left( \frac{\partial \theta_i}{\partial \beta_k} \right) + \sum \frac{y_i - b'(\theta_i)}{\phi_i} \frac{\partial^2 \theta_i}{\partial \beta_j \partial \beta_k}$$

- $E \left( -\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} \right) = \sum \frac{V(\mu_i)}{\phi_i} \frac{x_{ij}}{g'(\mu_i)V(\mu_i)} \frac{x_{ik}}{g'(\mu_i)V(\mu_i)} = \sum \frac{x_{ij}x_{ik}}{\phi_i \{g'(\mu_i)\}^2 V(\mu_i)}$



## ... maximum likelihood equation

- $\hat{\beta} = \beta + \mathbf{i}(\beta)^{-1} \mathbf{X}^T \mathbf{u}(\beta)$

$$\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} = \sum \frac{-b''(\theta_i)}{\phi_i} \left( \frac{\partial \theta_i}{\partial \beta_j} \right) \left( \frac{\partial \theta_i}{\partial \beta_k} \right) + \sum \frac{y_i - b'(\theta_i)}{\phi_i} \frac{\partial^2 \theta_i}{\partial \beta_j \partial \beta_k}$$

- $E \left( -\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} \right) = \sum \frac{V(\mu_i)}{\phi_i} \frac{x_{ij}}{g'(\mu_i)V(\mu_i)} \frac{x_{ik}}{g'(\mu_i)V(\mu_i)} = \sum \frac{x_{ij}x_{ik}}{\phi_i \{g'(\mu_i)\}^2 V(\mu_i)}$

$$\begin{aligned} \hat{\beta} &= \beta + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{u}(\beta) = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{W} \mathbf{X} \beta + \mathbf{X}^T \mathbf{u}(\beta) \} \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{W} (\mathbf{X} \beta + \mathbf{W}^{-1} \mathbf{u}(\beta)) \} \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \end{aligned}$$

## ... maximum likelihood equation

- $\hat{\beta} = \beta + \mathbf{i}(\beta)^{-1} \mathbf{X}^T \mathbf{u}(\beta)$

$$\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} = \sum \frac{-b''(\theta_i)}{\phi_i} \left( \frac{\partial \theta_i}{\partial \beta_j} \right) \left( \frac{\partial \theta_i}{\partial \beta_k} \right) + \sum \frac{y_i - b'(\theta_i)}{\phi_i} \frac{\partial^2 \theta_i}{\partial \beta_j \partial \beta_k}$$

- $E \left( -\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} \right) = \sum \frac{V(\mu_i)}{\phi_i} \frac{x_{ij}}{g'(\mu_i)V(\mu_i)} \frac{x_{ik}}{g'(\mu_i)V(\mu_i)} = \sum \frac{x_{ij}x_{ik}}{\phi_i \{g'(\mu_i)\}^2 V(\mu_i)}$

$$\begin{aligned} \hat{\beta} &= \beta + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{u}(\beta) = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{W} \mathbf{X} \beta + \mathbf{X}^T \mathbf{u}(\beta) \} \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{W} (\mathbf{X} \beta + \mathbf{W}^{-1} \mathbf{u}(\beta)) \} \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \end{aligned}$$

- does not involve  $\phi_i$

## ... maximum likelihood equation

- $\hat{\beta} = \beta + \mathbf{i}(\beta)^{-1} \mathbf{X}^T \mathbf{u}(\beta)$

$$\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} = \sum \frac{-b''(\theta_i)}{\phi_i} \left( \frac{\partial \theta_i}{\partial \beta_j} \right) \left( \frac{\partial \theta_i}{\partial \beta_k} \right) + \sum \frac{y_i - b'(\theta_i)}{\phi_i} \frac{\partial^2 \theta_i}{\partial \beta_j \partial \beta_k}$$

- $E \left( -\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} \right) = \sum \frac{V(\mu_i)}{\phi_i} \frac{x_{ij}}{g'(\mu_i)V(\mu_i)} \frac{x_{ik}}{g'(\mu_i)V(\mu_i)} = \sum \frac{x_{ij}x_{ik}}{\phi_i \{g'(\mu_i)\}^2 V(\mu_i)}$

$$\begin{aligned} \hat{\beta} &= \beta + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{u}(\beta) = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{W} \mathbf{X} \beta + \mathbf{X}^T \mathbf{u}(\beta) \} \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{W} (\mathbf{X} \beta + \mathbf{W}^{-1} \mathbf{u}(\beta)) \} \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \end{aligned}$$

- does not involve  $\phi_i$

- iteratively re-weighted least squares

$W, z$  both depend on  $\beta$

## ... maximum likelihood equation

- $\hat{\beta} = \beta + \mathbf{i}(\beta)^{-1} \mathbf{X}^T \mathbf{u}(\beta)$

$$\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} = \sum \frac{-b''(\theta_i)}{\phi_i} \left( \frac{\partial \theta_i}{\partial \beta_j} \right) \left( \frac{\partial \theta_i}{\partial \beta_k} \right) + \sum \frac{y_i - b'(\theta_i)}{\phi_i} \frac{\partial^2 \theta_i}{\partial \beta_j \partial \beta_k}$$

- $E \left( -\frac{\partial^2 \ell(\beta; \mathbf{y})}{\partial \beta_j \partial \beta_k} \right) = \sum \frac{V(\mu_i)}{\phi_i} \frac{x_{ij}}{g'(\mu_i)V(\mu_i)} \frac{x_{ik}}{g'(\mu_i)V(\mu_i)} = \sum \frac{x_{ij}x_{ik}}{\phi_i \{g'(\mu_i)\}^2 V(\mu_i)}$

$$\begin{aligned} \hat{\beta} &= \beta + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{u}(\beta) = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{W} \mathbf{X} \beta + \mathbf{X}^T \mathbf{u}(\beta) \} \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{W} (\mathbf{X} \beta + \mathbf{W}^{-1} \mathbf{u}(\beta)) \} \\ &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z} \end{aligned}$$

- does not involve  $\phi_i$

- iteratively re-weighted least squares

- **derived response**  $\mathbf{z} = \mathbf{X} \beta + \mathbf{W}^{-1} \mathbf{u}$

$\mathbf{W}, \mathbf{z}$  both depend on  $\beta$

linearized version of  $\mathbf{y}$