

2 WEEK 1 SEMESTER 2

9/1/16. APPLIED LINEAR STATISTICAL METHODS 2

Other name : Generalised linear models (GLM)

Introduction

Linear regression and how it can be generalised

Usually given a set of data:

$$\{y^{(i)}, x^{(i)}\}_{i=1, \dots, n}$$

↑ vector called the explanatory variable

↑  $y \in \mathbb{R}$  response variable

for each variable define  $y^{(i)}, x^{(i)}$  → variable vector

$$y^{(1)}, x^{(1)} \Rightarrow y_1, x_1$$

$$y^{(2)}, x^{(2)} \Rightarrow y_2, x_2$$

$$y^{(n)}, x^{(n)} \Rightarrow y_n, x_n$$

A variable and vector!

For each data point associate a variable/vector.

LR: linear relationship  $y_i = \beta^T x_i + \epsilon_i$   $\beta$  - set of parameters  
 $\epsilon$  - noise in equation

$$y_1 = \beta^T x_1 + \epsilon_1$$

$$y_n = \beta^T x_n + \epsilon_n$$

Error  $\epsilon_i \sim N(0, \sigma^2)$   $i \in 1, \dots, n$   
 $\epsilon_i \perp \epsilon_j$   $i \neq j$  (residuals independent)

First equation

$$y_i = \beta^T x_i + \epsilon_i \sim N(\mu, \sigma^2)$$

$(y_i)$  only one observation available  
 $(x_i)$  only one obj. available for  $y_i$

2

16 A

b E

D

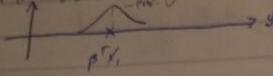
E

Question: What is  $p(y_i | \beta^T x_i)$ ?

$p$ : probability density function or distribution

$| \beta^T x_i$ : given  $\beta$  and  $x_i$  (known values)

$\Rightarrow$  What is the uncertainty associated with  $y_i$ ?



Will be a transformed normal dist with mean  $\beta^T x_i$ .

$$\begin{aligned} E[y_i] &= E[X_i^T \beta + \varepsilon_i] \\ &= E[X_i^T \beta] + E[\varepsilon_i] \quad \text{taken as given} \\ &= X_i^T \beta + 0 \end{aligned}$$

$$\text{Var}[y_i] = \sigma^2$$

$$\begin{aligned} p(y_i | \beta, x_i) &= N(y_i; \beta^T x_i, \sigma^2) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2}\right] \end{aligned}$$

(can redefine Linear Regression as  
redefine  $y_i$  is normally distributed such that:  $p(y_i | \beta, x_i)$ )

$$E[y_i] = \beta^T x_i \quad E[(y_i - \beta^T x_i)^2] = \text{Var}[y_i] = \sigma^2$$

## GLM

$p(y_i | \beta, x)$  is a member of ~~and~~ the exponential family of distributions.

SLR  $\rightarrow$  One element is Normal/Gaussian distribution  $y \in \mathbb{R}$

Logistic Regression  $\rightarrow$  Bernoulli distribution and outcome  $y \in \{0, 1\}$

Poisson Regression  $\rightarrow$  Binomial distribution generalization of Bernoulli  $0 \leq y \leq n$  integer number  $y \in \{0, \dots, n\}$

Poisson Regression  $\rightarrow$  Poisson Distribution discrete, can go to infinity  $y \in \mathbb{I}$

19/11/16 APPLIED LINEAR STATISTICAL METHODS II

3

↳ Exponential Distribution models  $y \in \mathbb{R}^+$  continuous-positive real number like time to failure  
 Weibull Distribution Survival analysis

Exponential Family of Distributions

$$(y|\theta) = \exp [a(y) b(\theta) + c(\theta) + d(y)]$$

function  $a$  of parameter  $y$  function  $c$  of  $\theta$

"Any distribution that you can write down in this form is a member of family"

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

$$\exp \left[ \frac{-|y-\mu|^2}{2\sigma^2} - \log(\sqrt{2\pi}\sigma) \right] \quad \theta \text{ associated with application of } y$$

$$\exp \left[ \underbrace{\frac{-y^2}{2\sigma^2} + \frac{2y\mu}{2\sigma^2} + \frac{\mu^2}{2\sigma^2}}_{a(y)} - \log(\sqrt{2\pi}\sigma) \right]$$

$$a(y) b(\theta)$$

$E[y]$  is defined w.r.t.  $\theta$ .

Exponential Distribution

$$\begin{aligned} p(y|\theta) &= \theta \exp(-\theta y) & y \in \mathbb{R}^+ \\ &= \exp(-\theta y + \log \theta) & \theta \in \mathbb{R}^+ \text{ and non-zero} \\ &\quad b(\theta) \quad a(y) \quad c(\theta) \quad f(y)=0 \end{aligned}$$

pdf must integrate to 1

$$\int_0^{+\infty} p(y|\theta) dy = 1$$

$$\begin{aligned} \int_0^{+\infty} \theta \exp(-\theta y) dy &= \left[ -\theta \exp(-\theta y) \right]_0^{+\infty} \\ &= -\exp(-\theta \cdot 0) + \exp(0) \\ &= 1 \end{aligned}$$

4

What is  $E[y] = \int_0^\infty y p(y|\theta) dy$  y times the function  
 $= \int_0^\infty y \cdot \lambda e^{-\lambda y} dy$   
 $E[y] = \lambda$   $\lambda$  cannot be zero

Weibull Distribution  $y \in \mathbb{R}^+$   $\theta \in \mathbb{R}^{++}$   $\lambda \in \mathbb{R}^{++}$   
 $p(y|\theta) = \lambda y^{\lambda-1} \exp(-\lambda y)$   
 when  $\lambda=1$  weibull = exponential dist

Poisson Distribution  $y \in \mathbb{N}$   $\theta \in \mathbb{R}^{++}$   
 $p(y|\theta) = \frac{\theta^y}{y!} \exp(-\theta)$

$$\begin{aligned} E[y] &= \sum_{y \in \mathbb{N}} y \frac{\theta^y}{y!} \exp(-\theta) = \theta \\ &= \exp(-\theta) \theta \sum_{y=1}^{\infty} \frac{\theta^y}{(y-1)!} \text{ Taylor expansion} \\ &= \exp(-\theta) \exp(\theta) \\ &= \theta \end{aligned}$$

$\theta$  is where information about explanatory variable will be inserted

$\theta$  is the parameter used to plug in information from explanatory variable  $x$   
 $\theta \rightarrow (x^T \beta)$  somehow associated.

Example: In normal distribution  $E[y] = \mu = \theta = x^T \beta$  (linear regression).

$\theta = g(x^T \beta)$  (Link function, identity function in this case)

$\mathbb{R} \rightarrow \mathbb{R}$  mapping to different domain.

will want to map to outcome domain like exp and log, must be invertible

$g$  and  $g^{-1}$  need to be available

6 APPLIED LI

GLM

1. Have a set of variables  $X_i$ 

2. Response variable

3. Model (or)

4.  $g: m$ 

5. Estimation

6. B

9/1/16 APPLIED LINEAR STATISTICAL METHODS 2

### GLM

1. Have a set of independent responses  $y_i$  with associated explanatory variable  $x_i$ :  $\{(x_i, y_i)\}^n_{i=1}$

2. Response  $y_i \sim p(y_i | \theta_i)$  member of the exponential family

3. Model (construction)  $\stackrel{\text{link function}}{g}(\mathbb{E}[y_i]) = x_i^\top \beta$   
 $\mathbb{E}[y_i] = g^{-1}(x_i^\top \beta)$

4.  $g$ : monotonic differentiable function (ensures inverse existence)

5. Estimate  $\beta$  by  $\hat{\beta} = \arg \max \text{likelihood}$  or  $\text{argmax posterior probability}$ .

6.  $\mathbb{E}[y] = g^{-1}(x^\top \beta)$

21/01/16 APPLIED LINEAR STATISTICAL METHODS 2

$y$  response variable  $\dim(y) = 1$  + dimension  
 $\rightarrow E[y] \rightarrow g(x|\beta)$  link function  
 $y \sim P_{y|x}(y|x)$   
 ↳ probability density function/distribution  
 continuous r.r. discrete r.v.

Example: In L.R. using N dist  
 $y \sim \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{(y-\bar{x}\beta)^2}{2\sigma^2}\right)$   $g$ : identity function

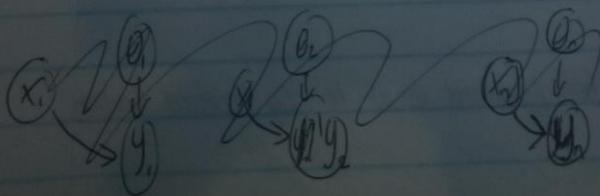
In general:  $y \sim P_{y|x}$   $\rightarrow P_{y|x,\beta}$   
 Some dist: weibull, exponential, poisson...  
 ↳ constrained model

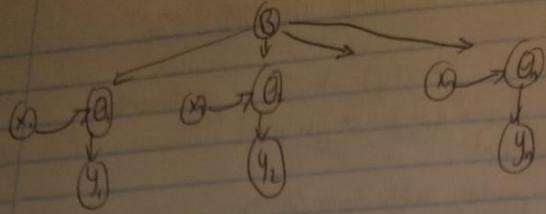
e.g. grade from hairs put in with parameters  $\beta$ .

Say Data  $n$  observations have mean in exm.  $\rightarrow$  Variable  
 Definition of variables  
 $x^{(1)} y^{(1)} \rightarrow y_1 | x_1 \rightarrow y_1 \sim P_{y|x}(y_1|x_1)$  Parameterized by  $\theta$   
 $x^{(2)} y^{(2)} \rightarrow y_2 | x_2 \rightarrow y_2 \sim P_{y|x}(y_2|x_2)$   
 $\vdots$   
 $x^{(n)} y^{(n)} \rightarrow y_n | x_n \rightarrow y_n \sim P_{y|x}(y_n|x_n)$   
 ↑  
 Value usually confused in note/test w/  $\int P_{y|x}(y|x_i) dy_i = 1$  mean value

$$\int P_{y|x}(y|x_i) dy_i = 1$$

Can't do this!





Control \$\theta\$ by \$\beta\$  
To estimate \$\theta\$ we only design a cost function

Likelihood \$\equiv\$ Cost function (correspond to joint dist of all values given params needed)

$$P(Y_1, Y_2, \dots | \theta_1, \theta_2, \dots, \theta_n)$$

ie your model is independent.

$$= \prod_{i=1}^n P(y_i | \theta_1, \theta_2, \dots, \theta_n)$$

independence of response variables given the \$\theta\$'s.

$= \prod_{i=1}^n P(y_i | \theta_i)$  only concerned with \$\theta\_i\$, other \$\theta\$ don't give information for \$y\_i\$.

$$\forall i: y_i \sim p_{y| \theta} (y_i | \theta_i)$$

Likelihood:  $P(Y_1, Y_2, \dots | \theta_1, \theta_2, \dots, \theta_n) = \prod_{i=1}^n p_{y| \theta} (y_i | \theta_i)$  Joint density function

$$\text{Total log likelihood: } \hat{\ell} = \sum_{i=1}^n \log(p_{y| \theta} (y_i | \theta_i))$$

center of exponential family of distributions

$$p_{y| \theta} (y_i | \theta) = \exp[a(y_i)b(\theta) + c(\theta) + d(y_i)]$$

$$\ell_{\text{loglik}} = \sum_{i=1}^n a(y_i)b(\theta_i) + c(\theta_i) + d(y_i)$$

each \$\theta\_i \rightarrow g^{-1}(x\_i^\top \beta)\$ link function

$$\beta = \arg \max \ell_{\text{loglik}}$$

find param that max the log likelihood function

2 21/01/16 APPLIED LINEAR STATISTICAL METHODS 2

3

GLMs : 1970's didn't have computational power

$$\text{LR: } y_i = \beta_0 + \beta x_i + \epsilon$$

$$= \begin{pmatrix} 1 & x_i \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \epsilon$$

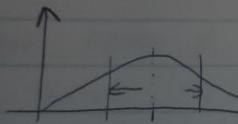
$x_i$  vector

$$y_i = \beta^T x + \epsilon$$

$\downarrow g$  used to get  $\theta_i$  using link function

$\theta_i$  can be to predict raw value

$$\beta \rightarrow P_{y|x,\beta} (y|x,\beta)$$



can get CI from it.

Only have one maximum for cost function (convex function)  
differentiate and equate to 0 for max value Only one global max solution

26/9/16 APPLIED LINEAR STATISTICAL METHODS 2.

### LOGISTIC REGRESSION (Ch 7 in Dobson book)

Usually:  $\{(x_i, y_i)\}_{i=1, \dots, n}$

For example:  $\{(x_i, y_i, n_i)\}_{i=1, \dots, n} [D(x_i) \text{ #Exp}(n_i), H(y_i)]$

Binomial Distribution:  $P(Y|\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$

$$y \in \{0, 1, \dots, n\} \quad \theta \in [0, 1]$$

$$\binom{n}{y} = \frac{n!}{(n-y)! y!} \sum_{y=0}^n \binom{n}{y} \theta^y (1-\theta)^{n-y}$$

Pascal's triangle

$$= (\theta + 1 - \theta)^n$$
$$= 1^n = 1$$

Binomial: member of exponential family of distributions?

$$\hookrightarrow \exp[a(y)b(\theta) + c(\theta) + d(y)]$$

If we can get binomial dist in this form? It prob 1/8 part of exponential family

$$\begin{aligned} P(Y|\theta) &= \exp[\ln \left[ \binom{n}{y} \theta^y (1-\theta)^{n-y} \right]] \\ &= \exp[\ln \left[ \binom{n}{y} + y \ln \theta + (n-y) \ln (1-\theta) \right]] \\ &= \exp[\ln \left[ \binom{n}{y} + \frac{y \ln \theta}{1-\theta} + n \ln (1-\theta) \right]] \\ &\quad \frac{\partial}{\partial y} \frac{\partial}{\partial \theta} \end{aligned}$$

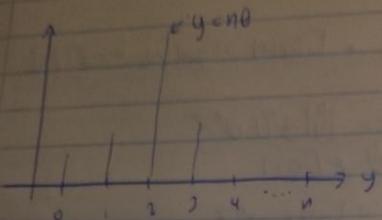
Value of  $\theta$  that maximise  $P(Y|\theta)$ ?

$$\begin{aligned} \frac{\partial P(Y|\theta)}{\partial \theta} &= \binom{n}{y} \left[ y \theta^{y-1} (1-\theta)^{n-y} - \theta^y (n-y) (1-\theta)^{n-y-1} \right] \\ &= \binom{n}{y} \theta^{y-1} (1-\theta)^{n-y-1} [y(1-\theta) - \theta(n-y)] \end{aligned}$$

$$\frac{d P(y|\theta)}{\theta} = 0 \quad y(1-\theta) - \theta(n-y) = 0$$

$$y - y\theta - \theta n + \theta y = 0$$

$$y = n\theta \quad \text{or} \quad \theta = y/n$$



$E[y]$ ? When  $y \sim P(y|\theta)$  (Binomial)

$$E[y] = \sum_{y=0}^n y \cdot P(y|\theta) = \sum_{y=0}^n y \cdot \binom{n}{y} \theta^y (1-\theta)^{n-y}$$

$$= \sum_{y=1}^n \underbrace{y \cdot \binom{n}{y} \theta^y (1-\theta)^{n-y}}_{\text{because when } y=0 \text{ the expression would have been 0.}}$$

$$y \binom{n}{y} = \frac{y n!}{(n-y)! y!} = \frac{y \cdot n \cdot (n-1)!}{(n-1-(y-1))! y \cdot (y-1)!} = n \binom{n-1}{y-1}$$

$$= n\theta \sum_{y=0}^{n-1} \binom{n-1}{y-1} \theta^{y+1} (1-\theta)^{(n-1)-(y-1)} = n\theta$$

$\underbrace{\qquad\qquad\qquad}_{\text{Bin}(\theta, n-1)}$

$$y_1 \sim \text{Bin}(\theta_1, n_1)$$

$$y_2 \sim \text{Bin}(\theta_2, n_2)$$

$$y_8 \sim \text{Bin}(\theta_8, n_8)$$

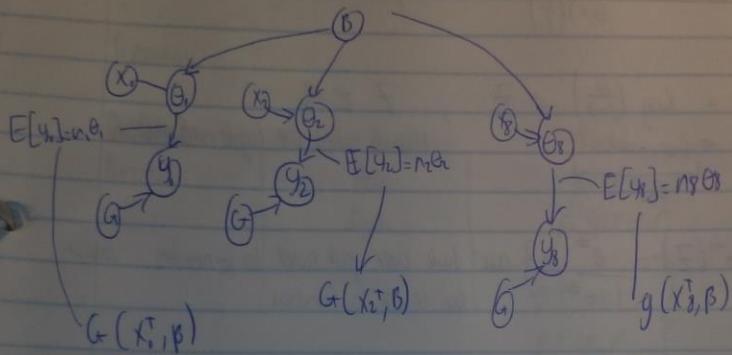
2

6/01/16 APPLIED LINEAR STATISTICAL METHODS 2

3

Saturated Model / Unbiased GLM

g-link function



$$L(\theta, \theta_1, \dots, \theta_4) = r(y_1, y_2, y_3, y_4 | \theta, \theta_1, \theta_2, \theta_3)$$

$$= \prod_{i=1}^4 r(y_i | \theta_i) \quad (\text{independence})$$

$$= \prod_{i=1}^4 \left( \frac{n_i}{y_i} \theta_i^{y_i} (1-\theta_i)^{n_i-y_i} \right)$$

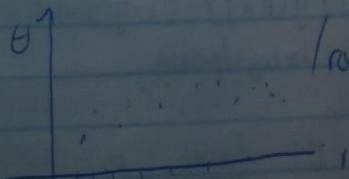
$$(\theta, \theta_1, \dots, \theta_4) = \underset{\theta, \theta_1, \dots, \theta_4}{\text{argmax}} L(\theta, \theta_1, \dots, \theta_4)$$

$$\theta_1 = y_1/n_1$$

$$\theta_2 = y_2/n_2$$

$$\theta_3 = y_3/n_3$$

Can plot \$\theta\_i\$ vs \$x\_i\$



(roughly)

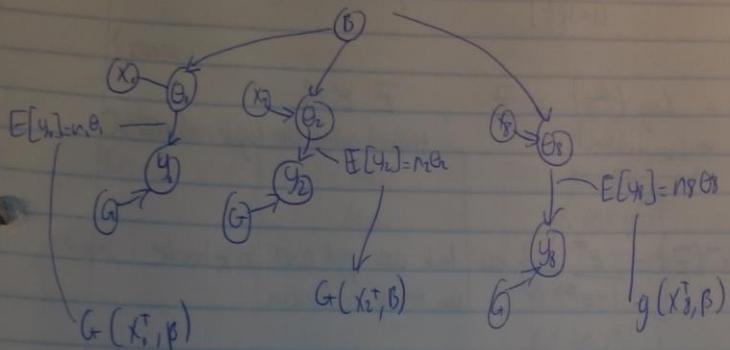
26/01/16

## APPLIED LINEAR STATISTICAL METHODS 2

3

Saturated Model / Unbiased GLM

g-linear model



$$\mathcal{L}(\theta_1, \theta_2, \dots, \theta_3) = P(y_1, y_2, y_3 | \theta_1, \theta_2, \theta_3)$$

$$= \prod_{i=1}^3 P(y_i | \theta_i) \quad (\text{independence})$$

$$= \prod_{i=1}^3 \left( \frac{n_i}{y_i} \theta_i^{y_i} (1-\theta_i)^{n_i-y_i} \right),$$

$$(\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_3) = \text{argmax } \mathcal{L}(\theta_1, \dots, \theta_3)$$

$$\hat{\theta}_1 = y_1/n_1$$

$$\hat{\theta}_2 = y_2/n_2$$

$$\hat{\theta}_3 = y_3/n_3$$

Can plot  $\theta_i$  vs  $x_i$  (roughly)

$$\theta \in [0, 1] \xrightarrow{g^{-1}} X, \beta \in \mathbb{R}$$

$\leftarrow g^{-1} \quad \text{or} \quad (1-x) \begin{pmatrix} \beta \\ 0 \end{pmatrix}$

Has to be invertible to transform both  
through the function to  $[0, 1]$

$$b(\theta) = \log \frac{\theta}{1-\theta} = z, \quad z \in \mathbb{R}$$

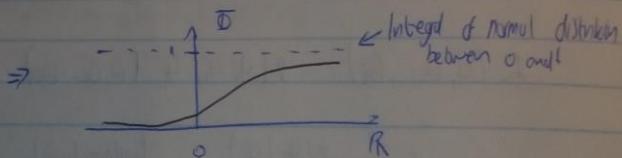
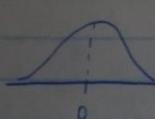
$b^{-1} : \mathbb{R} \rightarrow (-\infty, \infty)$

This is called the logit function

$$b^{-1}(z) = \frac{e^z}{1+e^z}$$

} nice! look over previous to remember how  
to inverse a function

Normal Distribution



The above is called the probit and it is positive and strictly increasing, therefore  
the linear case i.e.  $\beta \rightarrow D$

$$P(y| \beta x) = \binom{n}{y} g(\beta^T x)^y (1-g(\beta^T x))^{n-y}$$

$$\text{L } \beta = \prod_{i=1}^n P(y_i | \beta x_i)$$

$$= \prod_{i=1}^n \binom{n_i}{y_i} \left[ g(\beta^T x_i) \right]^{y_i} \left[ 1 - g(\beta^T x_i) \right]^{n_i - y_i}$$

$\hookrightarrow \text{Log-lik/ln/lnf/llk}$

Find  $\hat{\beta}$  such that  $\frac{\delta L(\beta)}{\delta \beta} = 0$ .

26/01/16 APPLIED LINEAR STATISTICAL METHODS 2

5

Note: You could add a  $\beta_2 x^2$  if you wanted to add another dimension to your model.

Summary:

Distribution	Link	Model
Binomial	Poisson	$\beta_0 + \beta_1 X$
	Logit	$\beta_0 + \beta_1 X + \beta_2 X^2$

$\beta_0$   
 $\beta_1 X$   
 $\beta_1 X + \beta_2 X^2$   
 $\beta_0 + \beta_2 X^2$

} in model creation

- Can use AIC & BIC to determine the "best" model
- Also must calculate the deviance to find out if the model is any good in the first place

28/01/16 APPLIED LINEAR STATISTICAL METHODS 2

GLM Case Study

$y_i$ : did patient survive? 1: NO 0: YES

$$\{ (x_i, y_i) \}_{i=1, \dots, 40} \quad n=40$$

$\uparrow$  age       $\uparrow$   $\varepsilon_{0,1}$

Response Variable  $y_i \sim \text{Binomial distribution}$

$$P(y_i|\theta) = \theta^{y_i} (1-\theta)^{1-y_i}$$

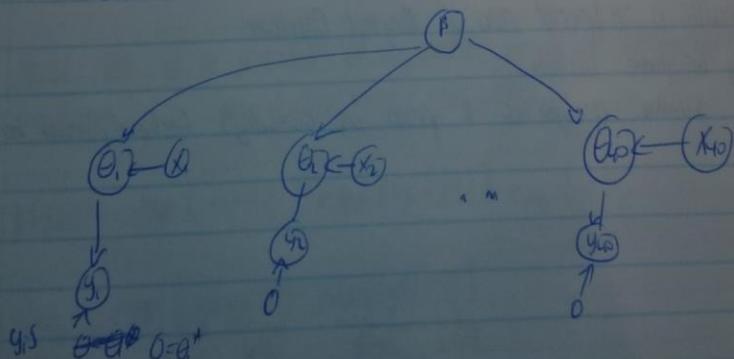
$$i=1 \quad P(y_1|\theta) = \binom{y_1}{\approx 0.13} \theta^{y_1} (1-\theta)^{1-y_1}$$

When there is only 1 choice (binary)  
The distribution used is called Binomial  
Bernoulli distribution

$$= \theta^{y_1} (1-\theta)^{1-y_1}$$

$$i=2 \quad P(y_2|\theta_2) = \theta_2^{y_2} (1-\theta_2)^{1-y_2}$$

$$i=40 \quad P(y_{40}|\theta_{40}) = \theta_{40}^{y_{40}} (1-\theta_{40})^{1-y_{40}}$$



$\Rightarrow$  likelihood function of saturated/unadjusted model with  $n=1$  (bernoulli)

$$E[y] = n\theta = \theta \quad n=1 (\text{bernoulli})$$

$$\frac{\delta P(y|\theta)}{\delta(\theta)} = 0 \Rightarrow \theta = y/n = y$$

$$l(\theta_0, \theta_1, \dots, \theta_n) = \prod_{i=1}^n \theta_i^{y_i} (1-\theta_i)^{1-y_i}$$

$\dim(\theta) = 4$

$$(\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_n) = \arg\max_{(\theta_0, \theta_1, \dots, \theta_n)} l(\theta_0, \theta_1, \dots, \theta_n)$$

$$\text{Link function: } b(\theta) - g(\theta) = \log(\frac{\theta}{1-\theta}) \quad [\text{Logit}]$$

$$\text{From } p(y|u) = \exp[u(y)b(u) + c(u) + d(y)]$$

GLM:

$$L(\beta) = \prod_{i=1}^n g^{-1}(\beta_0 + \beta_1 x_i) [1 - g^{-1}(\beta_0 + \beta_1 x_i)]^{1-y_i}$$

$\dim(\beta) = 2 \rightarrow$  calibrating the model

To test if a relationship exists in the model:

$$H_0: \beta_1 = 0 \quad \text{vs} \quad H_A: \beta_1 \neq 0$$

Summary:

- Bernoulli is a special case of Binomial Distribution
- Binary Response Variable
- Follow similar analysis of & proof as Tuesday's lecture (previous notes)

02/02/16 APPLIED LINEAR STATISTICAL METHODS 2

### Poisson Distribution - Poisson Regression

$$p(y|\lambda) = \frac{\lambda^y}{y!} e^{-\lambda} \quad y \in \mathbb{N} \quad \lambda > 0$$

Exercise:

1. Show this is a distribution. How? Show it is always positive.

$$p(y|\lambda) > 0 \rightarrow \text{positive function}$$

2. Show  $\sum_{y=0}^{\infty} p(y|\lambda) = 1$

$$= \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} e^{-\lambda} = e^{-\lambda} \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} \quad \text{we know it will integrate to 1!}$$

Aside: Taylor expansion  $e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!} \rightarrow$  this is what we have above.

$$\text{Taylor } e^x = \sum_{y=0}^{\infty} \frac{\lambda^y}{y!}$$

2. Show this is a member of exponential family of distributions.

In form of:  $\exp [a(y)b(\lambda) + c(\lambda) + d(y)]$

$$p(y|\lambda) = \exp [\log(\lambda^y)] \exp [-\log(y!)] \exp (-\lambda)$$

$$= \exp \left[ \underbrace{y \log(\lambda)}_{a(y)} + \underbrace{-\log(y!)}_{b(y)} + \underbrace{-\lambda}_{c(\lambda)} \right]$$

3. Compute the expectation of  $y$  and verify this is a function of  $\lambda$ .

$$\begin{aligned} \mathbb{E}[y] &= \int y p(y|\lambda) dy \quad \downarrow \text{change to sum because } y \text{ is a discrete r.v.} \\ &= \sum_{y=0}^{\infty} y p(y|\lambda) \end{aligned}$$

2.

$$\begin{aligned}
 & \sum_{y=0}^{\infty} \frac{y}{y!} \lambda^y \exp(-\lambda) \quad \text{take out } \exp(-\lambda) \text{ as does not depend on } y \\
 & = \exp(-\lambda) \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} \\
 & = \exp(-\lambda) \sum_{y=1}^{\infty} \frac{\lambda^y}{y(y-1)!} \quad \text{Start sum at } y=1 \text{ as at } y=0 \text{ sum } = 0 \\
 & = \exp(-\lambda + \lambda) \left[ \sum_{y=1}^{\infty} \frac{\lambda^{y-1}}{(y-1)!} \right] \quad (\text{Taylor expansion around } \lambda) \\
 & \sum_{y=1}^{\infty} \frac{\lambda^{y-1}}{(y-1)!} = \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} = \exp(\lambda) \\
 & = \exp(-\lambda)(\lambda)\exp(\lambda) \\
 & = \lambda = \text{expectation}
 \end{aligned}$$

4. Compute the maximum w.r.t  $\lambda$  of this distribution.

Find  $\lambda$  such that  $\frac{dP}{d\lambda} = 0$

$$\begin{aligned}
 \frac{dP}{d\lambda} &= \frac{1}{y!} \left[ y \lambda^{y-1} \exp(-\lambda) - \exp(-\lambda) \lambda y \right] \\
 &= \frac{\lambda^{y-1} \exp(-\lambda)}{y!} \left[ y - \lambda \right]
 \end{aligned}$$

equal to 0 when  $\lambda = y$  (from second bracket)

Sufficient soln of Max likelihood is  $\lambda = y$ .

$y$  is an integer,  $\lambda$  is in  $\mathbb{R}^+$

Log map  $\mathbb{R}$  onto  $\mathbb{R}^+$ , can use it as a link between our function and expectation.

$$E[y] = \lambda \in \mathbb{R} \quad \xrightarrow{\text{Log}(\lambda) = x^T \beta} x^T \beta$$

2.

1 only

02/02/16

## APPLIED LINEAR STATISTICAL METHODS 2

$$\mathbb{E}[y] = X \beta \quad \begin{matrix} \text{Log}(x) = x^T \beta \\ \lambda = \exp(x^T \beta) \end{matrix} \quad \begin{matrix} x^T \beta \\ \text{Link function} \end{matrix}$$

Are some links with poisson and binomial  
 Use poisson for number of occurrences of an event.  
 Sometimes have  $n$  poisson trials.

Binomial  $P(y|\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$   
 usually given observatory  $(y_i, x_i, n_i)$   $\mathbb{E}[y|\theta] = n\theta$

If we use log as link function  $\log[\mathbb{E}[y]] = \log(\lambda) = x^T \beta$   
 two possible models for r.v  $y$  (poisson or binomial).

Instead  $\log[\mathbb{E}[y]/n] = \log(\frac{\lambda}{n}) = x^T \beta$  like normalising (using size of group)  
 $\hookrightarrow \log(\lambda) = \log(n) + x^T \beta$  Divide expectation by size  $n$ .  
 equivalent of writing offset

Look at R code example  
 Some strategy at last week.

$n$  - is called exposure, "how many beds are you exposing to drug?"  
 $\log(n)$  is the offset.  $\Rightarrow$  only offset bc (intercept) in model

5pm 1.20 amnab doy a staff clinic

24.02.16

Beetle Data Distribution	Link	Explanatory Model	AIC
Binomial	Probit AIC: 40.52 Logit AIC: 41.41	$\hat{y} = \frac{\beta_0 + \beta_1 x_1}{1 + \exp(\beta_0 + \beta_1 x_1)}$ ; $P_i = \hat{P}_i^2$ ; $P_i = \hat{P}_i + \beta_2 x_2$	$-\ln \sum (\hat{P}_i)^{y_i} + 2m$ complexity penal with # of $\beta$ 's
Poisson	Log AIC: 57.72		

Insurance Dealer Example : See R code on website

Show that:

$$\lim_{n \rightarrow \infty} \text{Bin}(y; \theta, n) = \text{Pois}(y; \lambda)$$

$$\mathbb{E}[y] = n\theta \quad \text{Pois}(y; \lambda)$$

$$\lim_{n \rightarrow \infty} \binom{n}{y} \theta^y (1-\theta)^{n-y} \stackrel{?}{=} \frac{\lambda^y}{y!} \exp(-\lambda) \quad (1)$$

$$\binom{n}{y} \left(\frac{\lambda}{n}\right)^y \left(1 - \frac{\lambda}{n}\right)^{n-y} \xrightarrow{\text{expand}} = \frac{n^y}{A(n)} \frac{\lambda^y}{y!} \left(1 - \frac{\lambda}{n}\right)^y \left(1 - \frac{\lambda}{n}\right)^{-y}$$

$$\lim_{n \rightarrow \infty} A(n) \stackrel{?}{=} 1 \quad A(n) = \frac{\approx 1 (1-\lambda)}{\frac{n^y}{n} \dots \frac{(n-y+1)}{n}} \xrightarrow{n \rightarrow \infty} 1 \quad \text{dichit depends on } n \quad y \text{ goes to } 1$$

$$\lim_{n \rightarrow \infty} B(n) ? \quad \left(1 - \frac{\lambda}{n}\right)^n = \exp(-\lambda) \Rightarrow \left(1 - \frac{\lambda}{n}\right)^n = \sum_{k=0}^n \binom{n}{k} \left(\frac{-\lambda}{n}\right)^k$$

$$\begin{aligned} &\text{Put all back together} \\ &\left(1 - \frac{\lambda}{n}\right)^n / \exp(-\lambda) \xrightarrow{n \rightarrow \infty} = \text{poisson distribution} \end{aligned} \quad \begin{aligned} &\xrightarrow{\text{Taylor expansion}} \\ &= \frac{\frac{n!}{(n-k)! k!} \frac{-\lambda^k}{n^k}}{(n-k)! k!} \exp(-\lambda) \\ &= \frac{n!}{(n-k)! k!} \frac{-\lambda^k}{n^k} \quad \text{A(n) goes to 1} \end{aligned}$$

When I calc  $\Rightarrow$  Poisson

When  $n \rightarrow \infty$  calc  $\Rightarrow$  Poisson Distribution

07/02/16

## APPLIED LINEAR STATISTICAL METHODS 2

Looking at posterior odds

Example, villages and dirty water, how many are infected?

$n=10$  people total  
 $y=7$  infected

$H_0: \theta \leq 0.5$        $H_1: \theta > 0.5$   
not endemic      Infection is endemic.

Binomial:

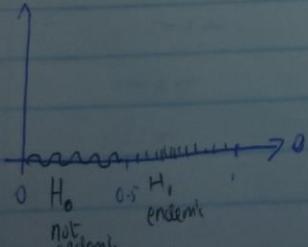
$$P(y|\theta) = \binom{10}{7} \theta^7 (1-\theta)^3$$

$$\hat{\theta}_b = y/n = 0.7 \text{ for Maximum Likelihood}$$

How to make 95% CI for this? of  $\theta$  value

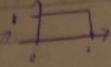
Posterior dist of  $\theta$  given  $y$ :

$$P(\theta|y) = \frac{p(y|\theta) \cdot p(\theta)}{p(y)} \propto \int p(y|\theta) p(\theta) d\theta \leftarrow \text{"normalizing constant"}$$



$$P = \int_{0.5}^1 p(\theta|y) d\theta \geq 0.95 \text{ then } H_1 \text{ is true}$$

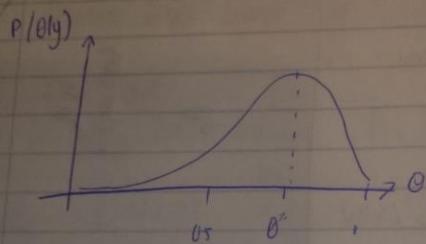
Distribution	$E[y]$	Max likelihood $\hat{\theta}^*$	$a(y)$	$b(y)$	$c(\theta)$	$d(y)$
Normal	$\frac{n\theta + y}{n+1}$	$\hat{\theta}_b = y$	$\hat{\theta}^* = \hat{\theta}_b$	$\hat{\theta}^* = \hat{\theta}_b$	$\hat{\theta}^* = \hat{\theta}_b$	$\hat{\theta}^* = \hat{\theta}_b$
Poisson	$\theta$	$\hat{\theta}_p = y$	$\hat{\theta}^* = \hat{\theta}_p$	$\hat{\theta}^* = \hat{\theta}_p$	$\hat{\theta}^* = \hat{\theta}_p$	$\hat{\theta}^* = \hat{\theta}_p$
Binomial	$\theta y$	$\hat{\theta}_b = y$	$\hat{\theta}^* = \hat{\theta}_b$	$\hat{\theta}^* = \hat{\theta}_b$	$\hat{\theta}^* = \hat{\theta}_b$	$\hat{\theta}^* = \hat{\theta}_b$
		$\hat{\theta}^* = \hat{\theta}_b$				

Belief ①:   $\theta$  can be anywhere between 0 and 1  
uniform probability for value of  $\theta$   
has to integrate to 1 on y axis

102/16

$$OP = \text{To find } \int_0^1 \left(\frac{\partial}{\partial \theta}\right) \theta^7 (1-\theta)^3 d\theta \quad P(\theta) = 1$$

$$\int_0^1 \theta^7 (1-\theta)^3 d\theta$$



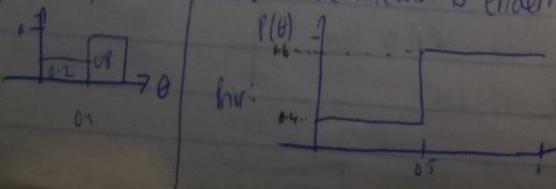
Checking that integral from  $0.5 \rightarrow 1$  is  $\geq 95\%$  to prove hypothesis.

$$\rightarrow \theta^7(1-\theta)^3 \Rightarrow -\frac{\theta^{11}}{11} + \frac{3\theta^{10}}{10} - \frac{3\theta^9}{9} + \frac{\theta^8}{8}$$

$$\begin{aligned} (1-\theta)^3 &= 1 - 3\theta + 3\theta^2 - \theta^3 \\ \theta^7(1-\theta)^3 &= \theta^7 - 3\theta^8 + 3\theta^9 - \theta^{10} \quad \text{integrate} \Rightarrow \\ &= \frac{\theta^8}{8} - \frac{3}{9}\theta^9 + \frac{3}{10}\theta^{10} - \frac{\theta^{11}}{11} \end{aligned}$$

$$\rightarrow \left[ \dots \right]_{0.5}^1 = 0.88 = \theta^*$$

(function): Belief ②: expert opinion: 80% sure that the infection is endemic



04/02/16

## APPLIED LINEAR STATISTICAL METHODS

3

Compute the prior again (the integral frusted)

$$P_{\text{expert}} = \frac{\int_{0.5}^1 \binom{10}{2} \theta^7 (1-\theta)^3 \cdot 1.6}{\int_{0.5}^1 \binom{10}{2} \theta^7 (1-\theta)^3 (0.4) + \int_{0.5}^1 \binom{10}{2} \theta^7 (1-\theta)^3 (1.6)} \leftarrow \text{split because no longer uniform}$$

Check if the is  $\geq 0.95$  to confirm H.

$$\theta^7 (1-\theta)^3 d\theta = \frac{-\theta^{11}}{11} + \frac{3\theta^{10}}{10} - \frac{3\theta^9}{9} + \frac{\theta^8}{8}$$

Results in 0.9691

Expert said money by not needing extra belts or larger n size to gain significant results.

11/02/16 APPLIED LINEAR STATISTICAL METHODS

### Choosing GLM Models

1 AIC for model selection

$$AIC = -2\log L_0 + 2M$$

Pick smallest AIC

$$\begin{aligned} \text{AIC Saturated} &= -2\log L_0 + 2\dim(\theta) \\ \text{GLM} &= -2\log L_p + 2\dim(\beta) \end{aligned}$$

!

2 Deviance: good Model?

$$\text{Deviance} = -2\log \left[ \frac{\hat{L}_B}{L_B} \right]$$

← Binomial Distribution  
← Binomial controlled by  $\beta^T X$

(Saturated Model)

For example, binomial distribution with  $\beta^T X$  and logit link  
 $\chi^2_{\text{dev}} \sim \text{chi-square}(n-p)$  Comment in R for chi-square test on Deviance

Example. Relation Between AIC and Deviance

$$\begin{aligned} \text{Deviance}_{\text{Binomial}} &= -2\log_{\theta} + 2\log_{\beta} \\ &= AIC_{\theta} - 2\dim(\theta) - AIC_{\beta} + 2\dim(\beta) \\ AIC_{\theta} &= -2\log_{\theta} + 2\dim(\theta) = AIC_{\theta} - AIC_{\beta} + 2(\dim(\beta) - \dim(\theta)) \end{aligned}$$

$$AIC_{\beta} = -2\log_{\beta} + 2\dim(\beta)$$

Chi-Square Distribution  
 $Z_i \sim N(0,1) \quad i=1, \dots, n$  independent i.i.d.

$$X = \sum_{i=1}^k Z_i^2 \quad X \sim \chi^2(k) \quad \text{prove this?}$$

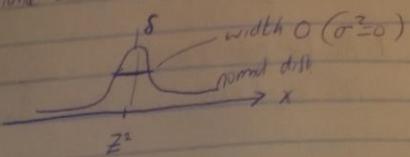
$$P_X(x) = \int_{-\infty}^{\infty} P_{XZ}(x, z) dz$$

When  $X = Z^2$   
 conditional distribution  
 $N(0,1)$

No uncertainty in  $x$  when given  $z$  ( $x$  is  $z^2$  reparam)

Dirac distribution has variance 0: no uncertainty

$$\underset{x|z}{p(x|z)} = \delta(x - z^2)$$



$$\int_{-\infty}^{\infty} \delta(x - z^2) dx = 1 \quad \text{property of dirac function}$$

$$= \int \delta(x - z^2) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) dz$$

Use Dirac property to rewrite dirac part

$$\Rightarrow \frac{1}{2\sqrt{\pi}} [\delta(\sqrt{x} + z) + \delta(\sqrt{x} - z)]$$

$$2|\sqrt{x}|$$

$$\Rightarrow \int \frac{1}{2|\sqrt{x}|} [\delta(\sqrt{x} + z) + \delta(\sqrt{x} - z)] \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{z^2}{2}\right] dz$$

$$= \frac{1}{2\sqrt{\pi}} \int_{-\infty}^{\infty} \delta(\sqrt{x} - z) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) dz + \int_{-\infty}^{\infty} \delta(\sqrt{x} + z) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) dz$$

Next property:  $\int_0^\infty f(t) \delta(t - T) dt = f(T)$

$$= \frac{1}{2|\sqrt{x}|} \left[ \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x}{2}\right) + \frac{1}{\sqrt{2\pi}} \exp\left(\frac{x}{2}\right) \right]$$

$$(\text{Chi-Square}): \text{for degrees } k \quad p_k(x) = \frac{1}{2^{k/2}} \frac{x^{k/2-1}}{\Gamma(k/2)} \exp\left[-\frac{x}{2}\right] \quad x \geq 0$$

$$\mathbb{E}[X] = k$$

Look for Deviance near value of  $k$  or on individual writing

ex:  $\chi^2$  for binomial ( $n$ ) / poisson ( $k$ ) etc

Re-check definition  
of chi-square dist  
may not be correct here

6/02/16 APPLIED LINEAR STATISTICAL METHODS

No lecture Thursday 8<sup>th</sup> feb

Smoking Dataset Example

Groups	Response Variable	"number per group"	$X_{age}$	$\hat{X}_{smoking}$	$\hat{P}_{smoking}$
i=1	$y_1 = 32$	$n_1 = 52407 = 1$	35-44	1=Yes	$P_1 = \frac{y_1}{n_1} = \frac{32}{52407}$
i=2	$y_2 = 104$	$n_2 = 4528 = 2$	45-54	1=Yes	
i=3	$y_3 = 206$	$n_3 = 28612 = 3$	55-64	1=Yes	$P_3 = \frac{y_3}{n_3} = \frac{206}{28612}$
i=4	$y_4 = 186$	$n_4 = 12663 = 4$	65-74	1=Yes	
i=5	$y_5 = 106$	$n_5 = 5317 = 5$	75-84	1=Yes	
i=6	$y_6 = 2$	$n_6 = 18790 = 1$	35-44	0=No	
i=7	$y_7 = 12$	$n_7 = 10473 = 2$	45-54	0=No	
i=8	$y_8 = 23$	$n_8 = 570 = 3$	55-64	0=No	
i=9	$y_9 = 20$	$n_9 = 2583 = 4$	65-74	0=No	
i=10	$y_{10} = 31$	$n_{10} = 1462 = 5$	75-84	0=No	

(in general counts size of group)  
or take middle value  
of category

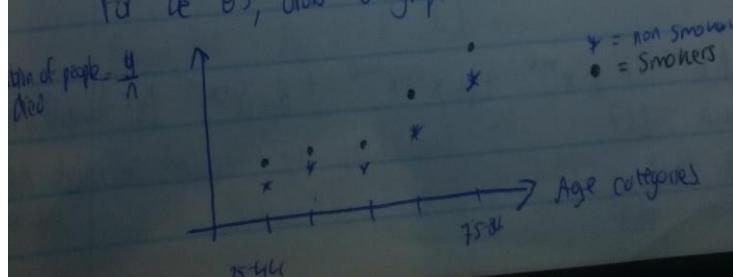
$N = 10$  (10 response variables)

Nature of explanatory variable:

- \* nominal (e.g. red, blue, green) (e.g. dead, alive or male, female) no order
- \* ordinal natural order between status of  $x$ .
- \* continuous e.g. weight, age, temperature

For the  $\hat{P}$ 's, draw a graph of all value to see difference

Separate into 2 graph



16/02/16 APPLIED LINEAR STATISTICAL METHODS

No lecture Thursday 18<sup>th</sup> feb

Smoking Dataset Example

Groups	Response Variable	(in general ranks size of group) "number per cent"	X <sup>age</sup>	X <sup>smoker</sup> 0=non 1=yes
i = 1	y <sub>1</sub> = 32	n <sub>1</sub> = 52407 = 1	35-44	1=Yes P <sub>1</sub> = $\frac{32}{52407}$
i = 2	y <sub>2</sub> = 104	n <sub>2</sub> = 428 = 2	45-54	1=Yes P <sub>2</sub> = $\frac{104}{52407}$
i = 3	y <sub>3</sub> = 206	n <sub>3</sub> = 2862 = 3	55-64	1=Yes P <sub>3</sub> = $\frac{206}{52407}$
i = 4	y <sub>4</sub> = 186	n <sub>4</sub> = 12663 = 4	65-74	1=Yes P <sub>4</sub> = $\frac{186}{52407}$
i = 5	y <sub>5</sub> = 106	n <sub>5</sub> = 5217 = 5	75-84	1=Yes P <sub>5</sub> = $\frac{106}{52407}$
i = 6	y <sub>6</sub> = 2	n <sub>6</sub> = 18790 = 1	35-44	0=No
i = 7	y <sub>7</sub> = 12	n <sub>7</sub> = 10673 = 2	45-54	0=No
i = 8	y <sub>8</sub> = 23	n <sub>8</sub> = 570 = 3	55-64	0=No
i = 9	y <sub>9</sub> = 20	n <sub>9</sub> = 2983 = 4	65-74	0=No
i = 10	y <sub>10</sub> = 31	n <sub>10</sub> = 1462 = 5	75-84	0=No

or take middle value  
of categories

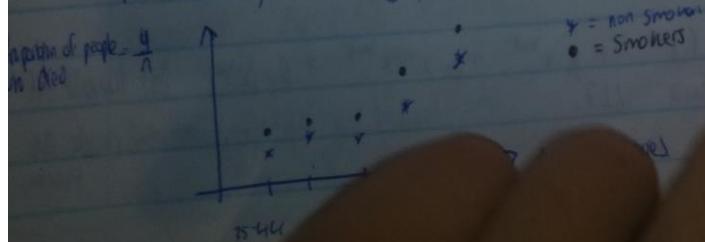
ordinal ↑ nominal ↑

N = 10 (10 response variables)

Nature of explanatory variable:

- \* nominal (e.g. red, blue, green) (e.g. dead, alive or male, female) no order
- \* ordinal natural order between states of X.
- \* continuous e.g. weight, age, temperature

For the X's, draw a graph of all relate to see difference



Separate into 2 group

GLM (one) to fit a curve between these points

Need to derive  $\beta^T x$

Could take  $age^2$  as an explanatory variable  $\Rightarrow 1, 2, 3, 4, 5 \rightarrow 1, 4, 9, 16, 25 \rightarrow$

Smoking is a binary variable,  $smoke^2$  has no effect

Multiply  $smoke$  by  $age$  (Product)  $\Rightarrow 1, 2, 3, 4, 5, 0, 0, 0, 0, 0$  mixed effect model

Multiply  $age^2 \times smoke$   $\Rightarrow 1, 4, 9, 16, 25, 0, 0, 0, 0, 0$  when two variables used

Point  $\rightarrow$  given two explanatory variables, can make new variables if relationship is not linear.

Normal, Binomial

$$\beta_0 + \beta_1 X^{age} + \beta_2 X^{smoke} + \beta_3 X^{age \cdot smoke} + \beta_4 (X^{age})^2 + \beta_5 (X^{smoke})^2$$

(original model from what we have defined in the table)

$$\begin{aligned} \text{When } smoke = 1: & (\beta_0 + \beta_2) + (\beta_1 + \beta_3)X^{age} + (\beta_4 + \beta_5)(X^{age})^2 \quad \text{"mixed model"} \\ \text{When } smoke = 0: & \beta_0 + \beta_1 X^{age} + (\beta_2 + \beta_4)(X^{age})^2 \end{aligned}$$

6 parameters here

Binomial or Poisson is the choice of distribution

Try all possible models and choose model with lowest AIC

Distr	Link	$x^T \beta$
Binomial	Logit	
Poisson	Log	

Fitted model:  $\beta_0 + \beta_1 X^{age} + \beta_2 X^{smoke} + \beta_3 X^{age \cdot smoke} + \beta_4 (X^{age})^2$  (using) lowest AIC

AIC: Poisson + Log link function  $66.7$  tried with  $Age^2$  removed  $\Rightarrow$  made AIC worse

Binomial + Logit  $66.63$

Binomial + probit =  $66.83$

16/02/16 APPLIED LINEAR STATISTICAL METHODS

3

When mixture effect term was removed, AIC increased

AIC: Poisson+Log : 667

Binomial+Logit : 668

Write the GLM  $g(\theta) = \beta^T x$

$$\text{Poisson+Log: } \log(\theta) = \hat{\beta}_0 + \hat{\beta}_1 x^{\text{age}} + \hat{\beta}_2 x^{\text{smoke}} + \hat{\beta}_3 x^{\text{exp}} x^{\text{smoke}} + \hat{\beta}_4 (x^{\text{exp}})^2$$

Binomial Logit  $\log\left(\frac{\theta}{1-\theta}\right) = \beta_0 + \beta_1 \dots$   $\beta$ 's will be different because of different Likelihood calculation

Deviances?  $-2 \log \left[ \frac{\text{Score}(L)}{\lambda \text{ GLM}} \right]$

Term  $-2 \log(L)$

1 For Saturated model  
- Poisson GLM  
- Binomial

2 For GLM Null Poisson+Log

3 For GLM Binomial+Logit.

$$\text{Poisson } (y | \lambda) = \lambda^y \exp(-\lambda) \quad (\text{10 groups } N=10)$$

$$L(\lambda_1, \lambda_2, \dots, \lambda_{10}) = \prod_{i=1}^{10} \frac{\lambda_i^{y_i}}{y_i!} \exp(-\lambda_i) \quad \lambda \text{ is expectation not proportion}$$

$y_i$  = dead people in each group i.

$$\log(L) \Rightarrow \sum_{i=1}^{10} y_i \log \lambda_i - \lambda_i - \log(y_i!)$$

$\hat{\lambda}_i = y_i$  Saturated model (Max likelihood Solution)

Deviance/AIC  $\log L(\lambda_1, \dots, \lambda_{10})$

$$\text{Binomial } p(y|\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$$

$$L(\theta_1, \theta_2, \dots, \theta_{10}) = \prod_{i=1}^{10} \binom{n_i}{y_i} \theta_i^{y_i} (1-\theta_i)^{n_i-y_i} \quad (\text{applied to each group})$$

02/16

ML Solution  $\hat{\theta} = y/m$  (saturated)Deviance/AIC:  $\text{Log}(L(\theta_0, \theta_{\text{obs}}))$ 

$$\text{Binomial} \\ L(\theta, \theta_{\text{obs}}) = \prod_{i=1}^m \binom{n_i}{y_i} \theta^{y_i} (1-\theta)^{n_i-y_i}$$

Saturated replace  $\theta$  with  $\theta = y/m$

Poisson

$$L(\lambda, \lambda_{\text{obs}}) = \prod_{i=1}^m \frac{\lambda^{y_i}}{y_i!} e^{-\lambda} [-\lambda]$$

Saturated  $\lambda_i = y_i$ 

$$\text{GLM}(y_i) \quad \tilde{\theta}_i = \exp(\hat{\beta}^T x_i) \quad \text{likelihood only dependent on } (\beta) \quad L(\beta)$$

(saturated)  $\tilde{x}_i = \exp(\hat{\beta}^T x_i)$

Differential related to age? Look at the coefficients  $\beta$ 's to determine  
 look at standard errors associated with the  $\beta$ 's.

"couldn't remove age explaining variable from model  $\Rightarrow$  factor,  
 variable is important."

Had to be used in the model.

Look at R code online

23/02/16

## APPLIED LINEAR STATISTICAL METHODS

Response Variable Member of Exponential Family

$y \in \{0, 1\}^n$  Bernoulli

$y \in \{0, 1, \dots, n\}$  Binomial

$y \in \text{integer}$  Poisson

$y \in \mathbb{R}$  Gaussian

$y \in \mathbb{R}^+$  exponential/weibull

positive real numbers (live time to failure)

$$\mathbb{E}[y] = \theta = g(x\beta)$$

Lukemia example.

Response variable  $y$  is time.

Is placebo time same as drug time?

Response variable for duration and any positive real number

### Survival Analysis

Modelling time until failure

Survival analysis is concerned with the statistical modelling of time to failure from a well defined origin (of time) or starting point

- time of hard drive to fail from the time it has been built/bought

- time of patient to die from the time disease has been diagnosed

#### 1. Distribution.

- times  $y$  are non-negative ( $y \in \mathbb{R}^+$ )

- distributions are skewed with long tails

- some subjects may survive beyond the study and their failure time is not observed.

In this case, the data is said to be "censored".

#### Exponential Distribution

$$p(y|\theta) = \theta \exp[-\theta y]$$

$$y \in \mathbb{R}^+ \quad \theta \in \mathbb{R}^{+*} \quad \mathbb{R}^+ \setminus \{0\}$$

3/02/16

Expectation  $E[y]^?$

$$\begin{aligned} E[y] &= \int_{-\infty}^{\infty} y p(y) dy \\ &= \int_0^{\infty} y \theta \exp[-\theta y] dy \\ &= \frac{1}{\theta} \text{ by integration by parts} \end{aligned}$$

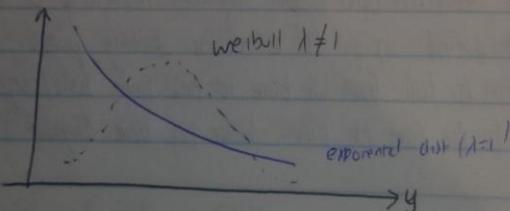
positive number

bell curve function to map  $\theta$  to positive real numbers  $\mathbb{R}^+$  onto  $\mathbb{R}$   
 $y$  should be exponential  $y^{-1}$  should be  $\log$

Weibull Distribution

$$p(y|\theta, \lambda) = \lambda \theta y^{\lambda-1} \exp[-\theta y^\lambda] \quad \begin{matrix} \leftarrow y \text{ power of } \lambda \\ \text{year } y \in \mathbb{R}^+ \quad \theta \in \mathbb{R}^+ \quad \lambda \in \mathbb{R}^+ \end{matrix}$$

Exponential Dist is a particular case of Weibull with  $\lambda=1$



$$\begin{aligned} E[y]? \quad E[y] &= \int_0^{\infty} y \lambda \theta y^{\lambda-1} \exp[-\theta y^\lambda] dy \\ &= \int_0^{\infty} \lambda \theta y^\lambda \exp[-\theta y^\lambda] dy \end{aligned}$$

\* EXAM PROB THIS

$$\text{Show } E[y] = \left(\frac{1}{\theta}\right)^{\frac{1}{\lambda}} \Gamma(1+\frac{1}{\lambda}) \quad \begin{matrix} \text{gamma} \\ \Gamma(\mu) = \int_0^{\infty} s^{\mu-1} \exp[-s] ds \end{matrix}$$

(A) Survivor and Hazard function

Probability of failure

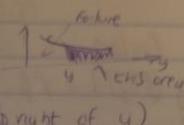
$$P(0 \leq t \leq y) = \int_0^y p(t|\theta) dt$$

p can be exponential or weibull area between 0 and y

23/02/16 APPLIED LINEAR STATISTICAL METHODS

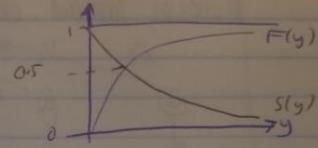
Probability of survival beyond time  $y$

$$S(y) = 1 - F(y) \quad (\text{area to right of } y)$$



B. Exercise  $F(y)$ ,  $S(y)$  & exponential distribution?

$$\begin{aligned} F(y) &= \int_0^y \theta \exp(-\theta t) dt \\ &= [-\exp(-\theta t)]_0^y \\ &= 1 - \exp(-\theta y) = \text{probability of failure.} \end{aligned}$$



$$S(y) = 1 - 1 + \exp(-\theta y) = \exp(-\theta y)$$

At 0.5 some probability of failure or survival  
where  $F(y) = 0.5$  or  $S(y) = 0.5$  called the median point of survival  
 $F(y) = S(y) = 0.5$

C. Hazard function

$$h(y) = \frac{p(y \leq t \leq y+dy)}{p(y \leq t)} \quad \begin{matrix} \text{chance of failure between time } y \text{ and } y+dy \\ \text{surviving beyond time } y \quad S(y) \end{matrix}$$

$$= p(y|t) dy$$

"chance of dying in next minute given a survival rate"

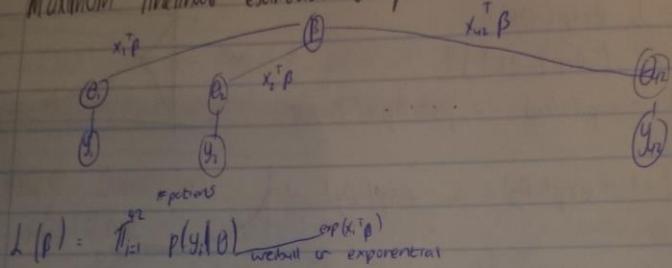
$$h(y) = \frac{p(y|t)}{S(y)} \quad \begin{matrix} \text{pdf} \\ \text{Weibull or exponential function} \end{matrix} \quad \begin{matrix} \text{hazard function} \\ \text{survival function} \end{matrix}$$

$$\text{Cumulative hazard function } H(y) = -\log [S(y)] \quad \text{primitve}$$

Link function  
 $g(\mathbb{E}[y]) = x^T \beta$

$\text{Log} \Leftrightarrow \theta + \exp(x^T \beta)$

Maximum likelihood estimation of  $\beta$ .



$$L(\beta) = \prod_{i=1}^{n+2} p(y_i | \theta) \quad \begin{matrix} \text{# patients} \\ \text{written or exponential} \end{matrix}$$

This is the definition of the likelihood when  $y_i$  is uncensored.

When  $y_i$  is censored, use  $s(y_i)$  instead of pdf  $p(y_i | \theta)$

$\delta_i$ : an indicator variable. 0 if  $y_i$  is censored, 1 if it's not.  $\delta_i \in \{0, 1\}$  censored

$$\hookrightarrow L(\beta) = \prod_{i=1}^{n+2} p(y_i | \theta)^{\delta_i} s(y_i)^{1-\delta_i}$$

-if you lose track of patient or  
 take either term depending on censored or not.  
 survives duration of study

Example exponential distribution.  $L(\theta_1, \dots, \theta_{n+2}) = \prod_{i=1}^{n+2} [\theta_i \exp(-\theta_i y_i)]^{\delta_i} [\exp(-\theta_i y_i)]^{1-\delta_i}$   
 (censored model)

$$L(\beta) = \prod_{i=1}^{n+2} [\exp(\beta^T x_i)]^{\delta_i}$$

Century variables given in example doct  
 GLM not defined for exponential or weibull  
 use "Survival" package

23/02/16 APPLIED LINEAR STATISTICAL METHODS

5

$\text{DF} = b_0 + b_1 = 2$  in example of tuberculosis  
Scale is  $1/\lambda$  value here

IS the drug better than the placebo effect?

look at

Median Survival time  $y_{\text{med}} = \frac{\log(2)}{\theta}$  for exponential dist

$$= \frac{\log(2)}{\exp(\hat{\beta}_0 + \hat{\beta}_1 x)}$$

$\xrightarrow{x=0 \text{ if placebo}}$   
 $x=1 \text{ drug}$

$$\frac{y_{\text{med}}|_{x=0}}{y_{\text{med}}|_{x=1}} = \frac{1}{\exp(-1.53)} = \exp(1.53) = \frac{1}{4} \Rightarrow \frac{y_{\text{med}}|_{x=0}}{4} = y_{\text{med}}|_{x=1}$$

$\frac{h|_{x=0}}{4} = h|_{x=1}$  4 times higher than treatment group. Probability is  
4 times higher with placebo rate

4 times more likely to fail.

R output gives  $\beta_1$ 's estimate on corresponding standard errors

$\pm 2 \cdot SE$

$2.32 \pm 2.04$  95% CI for  $\beta_1 \Rightarrow \beta_1$  is not 0.

~~Done~~

Same as in linear regression.  
relationship between time and drug exists

25/02/16 APPLIED LINEAR STATISTICAL

### Exponential Distribution

$$E[y] = \frac{1}{\theta}$$

$$\log [E[y]] = \log \left( \frac{1}{\theta} \right) = \beta_0 + \beta_1 x \quad \times \begin{cases} \text{positive / constant} \\ \text{map } \mathbb{R}^+ \text{ onto } \mathbb{R} \end{cases}$$

$$\log(\theta) = \beta_0 - \beta_1 x$$

$$\theta = \exp[-\beta_0 - \beta_1 x]$$

look at hazard function or median to see if drug is working

### Weibull

if true for weibull  $\rightarrow$  true for exponential

1. Show it is a distribution

$y$  positive,  $\theta$  positive,  $\exp$  is positive  $\rightarrow$  positive function  
must show that it integrates to 1

$$\int_0^\infty \lambda \theta y^{\lambda-1} \exp[-\theta y^{\lambda}] dy$$

$$= \left[ -\exp[-\theta y^{\lambda}] \right]_0^\infty \quad \text{at } 0 = 1$$

$$F(y) \quad \text{at } +\infty = 0 \Rightarrow \text{equal 1.}$$

$$\text{Probability of failure} = 1 - \exp[-\theta y^{\lambda}]$$

$$\text{Probability of survival} = 1 - F(y)$$

$$\text{Hazard function } h(y) = \frac{p(y|\lambda\theta)}{s(y)} = \lambda \theta y^{\lambda-1} \quad \text{for weibull (function of } y \text{)}$$

Exponential case is  $\lambda=1$   $h(y)$  is then  $\theta$ . (not a function of  $y$ ).  
independent of age "chance of survival" some at 100 or 1000 etc"

lack of memory ( $y$  important, many be a limitation)

2 Show it is a member of exponential distribution family

$$\begin{aligned} & \lambda y^{\lambda-1} \exp[-\lambda y^\lambda] \\ &= \exp[-\lambda y^\lambda + \log(\lambda) + \log(y^{\lambda-1})] \\ &= \exp[-\lambda y^\lambda + \log(\lambda) + \log(\lambda) + \log(y^{\lambda-1})] \\ & b(\alpha) \quad a(y) \quad c(\theta) \quad d(y) \quad \lambda \text{ is nuisance parameter} \end{aligned}$$

103 | 16

3 Expectation

$$E[y] = \int_0^\infty y \cdot \lambda y^{\lambda-1} \exp[-\lambda y^\lambda] dy$$

$$E[y] = \left(\frac{1}{\lambda}\right)^{\lambda} \Gamma\left(1+\frac{1}{\lambda}\right) \text{ with } \Gamma(u) = \int_0^\infty s^{u-1} \exp(-s) ds$$

$$\begin{aligned} \text{Substitution } u &= \lambda y^\lambda & du &= \lambda \lambda y^{\lambda-1} dy \\ &= \int_0^\infty y \cdot \lambda u \exp(-u) \end{aligned}$$

$$E[y] = \int_0^\infty \left(\frac{u}{\lambda}\right)^{\lambda} \exp(-u) \quad \text{look up proof online}$$

Inv function applied to expectation

Careful definition of weibull dist - in R package, be careful

have to look back to get it WARNING

$$\log \theta = -3.07 - 1.73x \quad (\text{the relationship}) \quad x \text{ is indicator for class}$$

The log re-scales the  $\beta$  values

$$2.48 + 1267 \quad (\text{R number})$$

$$\log\left(\frac{1}{\theta}\right)^{\lambda}$$

08/08/16 APPLIED LINEAR STATISTICAL METHODS

### MULTINOMIAL DISTRIBUTION

Extension to Binomial distribution

Binomial is a joint probability distribution

$$P(y_1, y_2, \dots, y_J | \theta_1, \theta_2, \dots, \theta_J, n) = \frac{n!}{y_1! y_2! \dots y_J!} \theta_1^{y_1} \theta_2^{y_2} \dots \theta_J^{y_J}$$

$$\vec{\theta} = [\theta_1, \theta_2, \dots, \theta_J]$$

$$p(\vec{y} | \vec{\theta}, n) = \frac{n!}{y_1! y_2! \dots y_J!}$$

J: number of categories

for example: J=3 yes, no, maybe

y<sub>1</sub>: # say yes

y<sub>2</sub>: # say maybe

y<sub>3</sub>: # say no

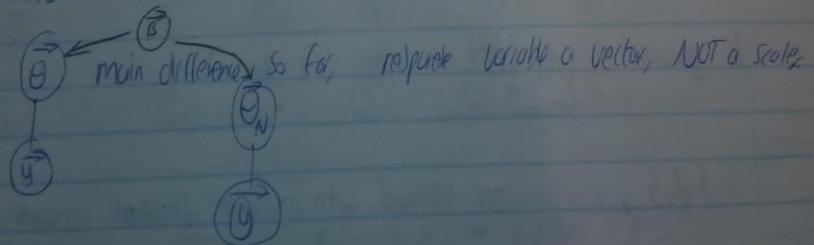
can see response variable as a vector.

constraint: n: total number of answers n = y<sub>1</sub> + y<sub>2</sub> + ... + y<sub>J</sub>

in example n = y<sub>1</sub> + y<sub>2</sub> + y<sub>3</sub> for J=3

constraint:  $\theta_1 + \theta_2 + \dots + \theta_J = 1$

$$E[y_1] = n\theta_1 \quad E[y_2] = n\theta_2 \dots$$



Exercise 1. J=2 Verify that multinomial is binomial distribution

$$P(y_1, y_2 | \theta_1, \theta_2, n) = \frac{n!}{y_1! y_2!} \theta_1^{y_1} \theta_2^{y_2} \quad n = y_1 + y_2 \text{ and } \theta_1 + \theta_2 = 1$$

$$P(y_1 | \theta_1, n) = \frac{n!}{y_1! (n-y_1)!} \theta_1^{y_1} (1-\theta_1)^{n-y_1} \quad (\text{Binomial distribution})$$

Exercise 2: is multinomial distribution a member of the exponential family  
of distributions? 16

$$\text{Exp: } p(y|t) = \exp [a(y) \cdot b(t) + c(t) + d(y)]$$

$t=2$ : yes, shown for binomial distribution

$t \geq 2$ : no

Even if multinomial distribution is not a member of the exponential family of distributions, we can control  $\vec{\theta}_y$  collected over  $N$  groups via a set of parameters  $\vec{\beta}$

### NOMINAL LOGISTIC REGRESSION

$$\log\left(\frac{\theta_j}{\theta_1}\right) = x^T \beta_j \quad \beta \text{ depends on category } j. \quad \begin{matrix} g=1 & \text{yes} \\ g=2 & \text{maybe} \\ g=3 & \text{no} \end{matrix}$$

$$\begin{cases} \log\left(\frac{\theta_2}{\theta_1}\right) = x^T \beta_2 & \text{associated with "maybe" } t=2 \\ \log\left(\frac{\theta_3}{\theta_1}\right) = x^T \beta_3 & \rightarrow \text{using } \theta_1 \text{ as a reference} \\ \dots & \beta \text{'s associated with "no" } t=3 \end{cases}$$

$$\text{In general } \log\left(\frac{\theta_j}{\theta_1}\right) = x^T \beta_j \quad j = 2, \dots, J$$

↑ reference category

$\{\beta_j\}_{j=2 \dots J}$  are estimated with maximum likelihood approach

$$\hat{\theta}_j = \frac{\exp(x^T \beta_j)}{1 + \sum_{j=2}^J \exp(x^T \beta_j)} \quad \forall j = 2, \dots, J$$

↳ SOFTMAX FUNCTION  
converting output to probability

$$\hat{\theta}_j = \frac{\exp(x^T \beta_j)}{\hat{\theta}_1 + \hat{\theta}_2 + \dots + \hat{\theta}_J} \quad (\text{taking exp of both sides})$$

08/03/16 APPLIED LINEAR STATISTICAL METHODS

Can rewrite as:

$$\beta_0 + \beta_1 \exp(x_1^T \beta) + \dots + \beta_k \exp(x_k^T \beta) = 1$$

$$\beta_j = \frac{1}{1 + \exp(x_1^T \beta) + \dots + \exp(x_k^T \beta)}$$

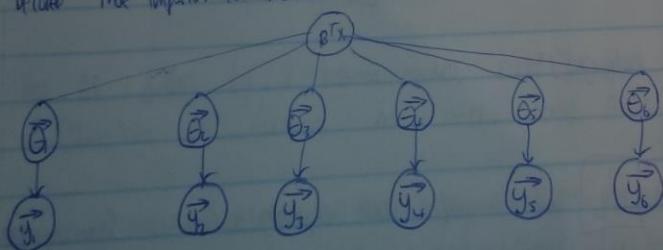
### EXAMPLE ON WEBSITE

Groups.

	$y_1$	$y_2$	$y_3$	$n$	GENDER	AGE or Age2
[Women 18-23]	i=1	26.2	12.5	$\frac{1}{n} = \frac{1}{45}$	$F=0$	0 0 0
[Women 24-40]	i=2	19.0	12.5	$\frac{1}{n} = \frac{1}{45}$	$F=0$	1 1 0
[Women >40]	i=3	15.5	14.4	$\frac{1}{n} = \frac{1}{60}$	$F=0$	2 0 1
[Men 18-23]	i=4	40.4	17.9	$\frac{1}{n} = \frac{1}{65}$	$M=1$	0 0 0
[Men 24-40]	i=5	17.7	15.5	$\frac{1}{n} = \frac{1}{44}$	$M=1$	1 1 0
[Men >40]	i=6	18.5	15.5	$\frac{1}{n} = \frac{1}{41}$	$M=1$	2 0 1
Saturated Model		$\theta_0$	$\theta_1$	$\theta_2$		

Picture: More important for women within older

indexed by  $i$



$$\text{Model 1: } \log\left(\frac{\theta_i}{\theta_0}\right) = \beta_{0j} + \beta_{1j} \text{Sex} + \beta_{2j} \text{Age}_1 + \beta_{3j} \text{Age}_2$$

$$\text{Model 2: } \log\left(\frac{\theta_i}{\theta_0}\right) = \beta_{0j} + \beta_{1j} \text{Sex} + \beta_{2j} \text{Age}$$

Model 4 has 2 less  $\beta$ 's hence lower AIC.

Model 2 has reduced complexity, not entirely a better fit.

4

3/16

What are  $\beta$ 's of Model 1?

use Softmax function

	$i=1$	$i=2$	$i=3$	$i=4$	$i=5$	$i=6$	$\beta_3$
				*	*		

$$\theta_j = \frac{\exp(x^T \beta_j)}{1 + \sum_{j=2}^6 \exp(x^T \beta_j)}$$

SOFTMAX FUNCTION

$$\hat{\theta}_{21} = \frac{\exp(-0.59 - 0.38\text{Sex} + 1.13\text{Age1} + 1.59\text{Age2})}{1 + \exp(-0.59 - 0.38\text{Sex} + 1.13\text{Age1} + 1.59\text{Age2}) + \exp(-1.04 - 0.81\text{Sex} + 1.48\text{Age1} + 2.92\text{Age2})}$$

(can do same for  $\theta_3$  and  $\theta_4$  is  $[-\hat{\theta}_2 - \hat{\theta}_3]$ )

group "i" category "j"

Should be  $\hat{\theta}_{2i}$  and i's under each explaining var i.e. Sex; Age;

Odd ratio:  $\frac{\hat{\theta}_{21}}{\hat{\theta}_{11}} = 1.5$   $\rightarrow$  year woman 18-23

$\frac{\hat{\theta}_{24}}{\hat{\theta}_{21}}$   $\rightarrow$  men 18-23

men important/not important

Odd RATIO assuming softmax function

Ratio of women if important over not important is bigger than that of men for age 18-23.

10/03/16

## APPLIED LINEAR STATISTICAL METHODS

### MULTINOMIAL DISTRIBUTION

Instead of one value of response, you have a vector of response ( $\vec{y}$ )

Observe a number of groups with vector responses, number of people in each group

$$\xi(\vec{y}_i, x_i, n_i) \quad i=1 \dots N \text{ (number of groups)}$$

$\downarrow$  number of people voting in group  $i$ .  
explanatory variable associated with group  $i$ .

J categories, J-1 degrees of freedom.

number of terms by J-1

i.e. intercept, sex, age =  $3 \times 2 = 6$  d.f. in example.

d.f. of saturated model 2 d.f. for each group (3 options per group)

6 groups = 12 d.f. for saturated model

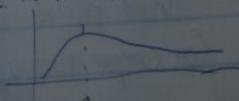
If AIC differs by 2 between models, could be from an extra variable  
(2m) in AIC formula, m is m+1 hence extra 2.

Comparing  $m_1$  and  $m_2$ ,  $m_1$  has 8 parameters,  $m_2$  has 6 parameters.

The difference between log likelihood  $-2[\text{Log } L(m_1) - \text{Log } L(m_2)] \approx 2$

Are the models that different?

$$= 2 \sim \chi^2(\text{d.f.} 2)$$



2 will be in 95% CI of chi-squared.

Can now say  $m_1 \approx m_2$  basically equivalent.

2

## ODD RATIO

$$OR = \frac{\theta_{j+2,i+1}}{\theta_{j+1,i+1}}$$

 $i = \text{group}$   $j = \text{catechase}$ 

- relative group: women ( $\theta_{j+2}$ )  
 ratio for group 4 men ( $\theta_{j+1}$ )  
 $\begin{cases} \text{Sex}=0 \\ \text{Sex}=1 \end{cases}$   
 $\begin{cases} \text{Age}=0 \\ \text{Age}=1 \end{cases}$

Category 2 =  $y_2$        $j=1$  = reference category

Only explanatory variable changing between groups will be seen.

OR =  $\exp[+0.38 \cdot \beta_0]$   $\approx 1.5$  indicates the option is more important to women than to men  
 (and logit link function)

$$\log \left[ \frac{\theta_{j+2,i+1}}{\theta_{j+1,i+1}} \right] = -0.59 - 0.38 \text{Sex}_i + 1.12 \text{Age}_i + 1.58 \text{Age}_2^2$$

Compute for  $i=1$  and  $i=4$  and compute ratioonly thing that will change is  $\text{Age}_i$  to  $-0.38 \cdot \text{Age}_i$ 

How to up standard errors to ensure odds ratio outcome is correct?

OR =  $\exp[+0.38 \pm 2 \cdot 0.38]$  0.38 standard error or Jex in  $y_2$

2SE of Sex (Sex was being tested with J)

$$CI \left[ (\exp[-0.38 - 2 \cdot 0.38]), (\exp[-0.38 + 2 \cdot 0.38]) \right] = [0.8, 2.08] \quad 95\% CI$$

bc OR related to only one explanatory variable on odds RATIO

1 is within the CI  $\Rightarrow$  men + women may be having the same.

Compare group 3 to 6 because age2 or age1 will change will get new results  
 Should provide some result of previous test. end up with some p

(on compare category 3 v category 1 (use  $y_3$  test p)). Will be outside that interval

"not important" chosen as reference

5/3/16

## APPLIED LINEAR STATISTICAL METHODS

Remark using Multinomial Distribution for nominal logistic Regression

$$\text{Logit}(\theta_j) \quad \log \left[ \frac{\theta_j}{\theta_1} \right] = \mathbf{x}^T \boldsymbol{\beta}_j \quad j=2, \dots, J \quad J = \# \text{ of categories (yes/no, male/female)}$$

softmax function

why log function?

$$\text{When } J=2: \log \left[ \frac{\theta_2}{\theta_1} \right] = \mathbf{x}^T \boldsymbol{\beta}_2 \quad \text{only equation we have (binomial dist)}$$
$$\theta_1 + \theta_2 = 1$$

Alligators Example Handout

A Suited because it allows for multiple response categories. In this case 5 different categories

B Three piece of information is enough to deduce which of four categories it belongs to. If  $L_1, L_2$  and  $L_3$  are all 0, then mean, it is Lake 4. Group

$$C \quad \log \left[ \frac{\theta_j}{\theta_1} \right] = \text{intensity} \quad j=2, 3, 4, 5 \\ \theta_0 + \beta_1 \text{Sex} + \beta_2 \text{Size} + \beta_3 \text{Lake}_1 + \beta_4 \text{Lake}_2 + \beta_5 \text{Lake}_3$$

$$D. \text{ for } j=2, 3, 4, 5 \\ \text{maximize } \log \left[ \frac{\theta_j}{\theta_1} \right] = \theta_0 + \beta_1 \text{Sex} + \beta_2 \text{Size} + \beta_3 \text{Lake}_1 + \beta_4 \text{Lake}_2 + \beta_5 \text{Lake}_3$$

$$E. \text{ AIC: Akaike information criterion} = -2 \log(L) + 2m$$

$L$  is likelihood function.  $M$  = number of parameters estimated  
It is a measure of model fit, rewarding model accuracy  
and penalizing extra parameters  
The lower the AIC, the better.

2

(03)16

C

$$\theta_j = \frac{\exp(\beta_j^T X)}{1 + \sum_{k=2}^J \exp(\beta_k^T X)} \quad j=2, \dots, J \quad \text{softmax formula}$$

$$\theta_1 = \frac{1}{1 + \sum_{k=2}^J \exp(\beta_k^T X)}$$

F:  $\theta_2 = 0.17 - 0.46\text{Sex} - 1.38\text{Size} - 1.78\text{Lchr}_1 + 0.41\text{Lchr}_2 + 1.16\text{Lchr}_3$   
 $\theta_3 = -3.42 - 0.63\text{Sex} + 0.56\text{Size} + 1.13\text{Lchr}_1 + 2.53\text{Lchr}_2 + 3.06\text{Lchr}_3$   
 $\theta_4 = -2.43 - 0.61\text{Sex} + 0.73\text{Size} + 0.58\text{Lchr}_1 - 0.55\text{Lchr}_2 + 1.24\text{Lchr}_3$   
 $\theta_5 = -1.43 - 0.25\text{Sex} - 0.3\text{Size} + 0.77\text{Lchr}_1 + 0.03\text{Lchr}_2 + 1.50\text{Lchr}_3$

G.

## SOLUTIONS

A. Multinomial distribution is an extension of the binomial distribution where the number of categories possible for the response is larger than 2. Here the number of categories is  $J=5$ , corresponding to 5 food choices available to all year 11 students.

$$F(Y_1, Y_2, Y_3, Y_4, Y_5 | \theta_1, \theta_2, \theta_3, \theta_4, \theta_5, n)$$

$$= \frac{n!}{Y_1! Y_2! Y_3! Y_4! Y_5!} \theta_1^{Y_1} \theta_2^{Y_2} \theta_3^{Y_3} \theta_4^{Y_4} \theta_5^{Y_5} \text{ with } Y_1 + Y_2 + Y_3 + Y_4 + Y_5 = n \\ \theta_1 + \theta_2 + \theta_3 + \theta_4 + \theta_5 = 1$$

B. Hancock:  $(l_1, l_2, l_3) = (1, 0, 0)$  4 different triplets to choose

Oklahoma: " " " = (0, 1, 0) 4! / 2! = 12

Trafford: " " " = (0, 0, 1)

Greater: " " " = (0, 0, 0)

2

3

15/03/16 APPLIED LINEAR STATISTICAL METHODS

$$C \text{ Log } \left[ \frac{\theta_j}{\theta_i} \right] = \beta_0^{\text{intercept}} + \beta_1^{\text{sex}} \text{ sex} + \beta_2^{\text{size}} \text{ size} + \beta_3^{\text{L1}} L_1 + \beta_4^{\text{L2}} L_2 + \beta_5^{\text{L3}} L_3$$

$$\forall j=2, \dots, 5.$$

$$\theta_j = \frac{\exp(\beta_j^T x)}{1 + \sum_{i=0}^{j-1} \exp(\beta_i^T x)} \quad \forall j=2, \dots, 5 \text{ with } x = [1, \text{sex}, \text{size}, L_1, L_2, L_3]$$

$$\beta_j = [\beta_0^{\text{intercept}}, \beta_1^{\text{sex}}, \beta_2^{\text{size}}, \beta_3^{\text{L1}}, \beta_4^{\text{L2}}, \beta_5^{\text{L3}}]$$

$$\theta_i = \frac{1}{1 + \sum_{j=0}^{i-1} \exp(\beta_j^T x)}$$

$$D \quad L = \prod_{i=1}^{16} P(Y_{1i}, Y_{2i}, Y_{3i}, Y_{4i}, Y_{5i} | \theta_{1i}, \theta_{2i}, \theta_{3i}, \theta_{4i}, \theta_{5i}, n_i)$$

$\sum n_i = Y_{1i} + Y_{2i} + Y_{3i} + Y_{4i} + Y_{5i}$  group sizes can be different  
Saturated likelihood of Scovard model

$$\theta_{1i} = \frac{1}{1 + \sum_{j=2}^5 \exp(\beta_j^T x_i)} \quad \leftarrow \text{used to calculate fitted values in table 8}$$

$$\theta_{ji} = \frac{\exp(\beta_j^T x_i)}{1 + \sum_{j=2}^5 \exp(\beta_j^T x_i)}$$

$$L = \prod_{i=1}^{16} \frac{n_i}{Y_{1i}! Y_{2i}! Y_{3i}! Y_{4i}! Y_{5i}!} \theta_{1i}^{Y_{1i}} \theta_{2i}^{Y_{2i}} \theta_{3i}^{Y_{3i}} \theta_{4i}^{Y_{4i}} \theta_{5i}^{Y_{5i}} \quad d.f = 4 \times 16$$

$$L(\theta_{1i}, \theta_{2i}, \theta_{3i}, \theta_{4i}, \theta_{5i}) = \sum_{i=1}^{16} h(\theta_{1i}) \quad d.f = 6 \times 4$$

E Criterion to use to select best model. Can we to drop parameter or not.

4

33/16

$$\begin{aligned}
 F \beta_1^T x &= 0.17 - 0.465x_1 - 1.35x_2 - 1.78x_3 + 0.91x_4 + 1.16x_5 \\
 \beta_2^T x &= -7.42 - 0.635x_1 + 0.565x_2 + 1.13x_3 + 2.53x_4 + 3.06x_5 \\
 \beta_3^T x &= -2.63 - 0.615x_1 + 0.735x_2 + 0.58x_3 - 0.55x_4 + 1.24x_5 \\
 \beta_4^T x &= -1.43 - 0.215x_1 - 0.245x_2 + 0.77x_3 + 0.03x_4 + 1.56x_5
 \end{aligned}$$

G Use table of fitted values

	Hancock	vs	Okamoto	Using Table 8 odd ratio
male small	i=1 $\frac{0.05}{0.6}$ but $\frac{0.08}{0.07}$	i=5 $\frac{0.08}{0.30}$ but $\frac{0.02}{0.02}$	3.08	
male big	i=2 $\frac{0.11}{0.12}$	i=6 $\frac{0.028}{0.48}$	3.08	
female small	i=3 $\frac{0.07}{0.071}$	i=7 $\frac{0.01}{0.21}$	3.08	
female Big	i=4 $\frac{0.17}{0.02}$	i=8 $\frac{0.058}{0.36}$	3.08	

$$\text{OR} = \exp \left( \beta_{j=4}^{L_1} - \beta_{j=4}^{L_2} \right) = 1.3$$

$$\begin{aligned}
 H \text{ (confident?)} \quad SE \quad \beta_{j=4}^{L_1} &= 0.79 \quad \beta_{j=4}^{L_2} = 1.20 \\
 \text{tale sum} \quad SE \left[ \beta_{j=4}^{L_1} - \beta_{j=4}^{L_2} \right]^2 &= \left[ E \left( \beta_{j=4}^{L_1} \right)^2 + E \left( \beta_{j=4}^{L_2} \right)^2 \right] \text{ independence assumed}
 \end{aligned}$$

$$SE = \sqrt{0.79^2 + 1.20^2} = 1.45$$

$1.13 \pm 1.45 \times 2 = 95\% \text{ CI}$   $0$  is in interval. Not confident.  
 expl=1, see if o chose that OR con or will be!  
 can't say for sure "algebraic ... in Okamoto than Hancock"

22/03/16 APPLIED LINEAR STATISTICAL METHODS

### INFERENCE

Algorithms to estimate  $\beta$ 's

Usually  $\hat{\theta} = \arg \max_{\theta} \text{lik}(\theta)$   
 $= \arg \max_{\theta} \log(\text{lik}(\theta))$  usually differentiable this function and equate to 0.

### Normal Distribution

link  $g$ : identity function [mapping  $R$  onto  $R$ ]  
 $p(y|\beta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[\frac{-(y-\beta)^2}{2\sigma^2}\right]$

$y_1, \dots, y_n$  :  $N$  responses

$x_1, \dots, x_n$

$x_i^T \beta$  (is constrained by error)  
 $\beta$  (constraining the mean)

likelihood wrt some parameter  $\beta$ :  $\text{lik}(\beta) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left[\frac{-(y_i-\beta)^2}{2\sigma^2}\right]$

$$\log(\text{lik}(\beta)) = -N \log(\sqrt{2\pi}\sigma) - \sum_{i=1}^N (y_i - x_i^T \beta)^2 / 2\sigma^2$$

$$SSE = \sum_{i=1}^N (y_i - x_i^T \beta)^2$$

$$\beta \rightarrow \frac{d(SSE)}{d(\beta)} = 0 \quad (\text{how we get our } \beta's)$$

$$SSE = (\|\vec{y} - X\beta\|)^2 \quad (\text{norm})^2$$

$$\vec{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}$$

$$\text{When we compute derivative of } SSE \quad \frac{d(SSE)}{d(\beta)} = \frac{d}{d(\beta)} (\vec{y} - X\beta)^T (\vec{y} - X\beta)$$

$$= \frac{d}{d\beta} [\vec{y}^T \vec{y} - \vec{y}^T X\beta - (X\beta)^T \vec{y} + (X\beta)^T X\beta] \quad \text{equal to zero}$$

$$= 0 + (\vec{y}^T X)^T - \vec{y}^T + 2 X^T X \beta = 0$$

$$\beta = (X^T X)^{-1} X^T \vec{y}$$

1. Differentiate log likelihood  $\log l(\theta)$

2. Set  $\nabla \log l(\theta) = 0$

3. Solving system of equations defined by  $\nabla \log l(\theta) = 0$

Newton-Raphson

looking at  $\log l(\theta)$  near  $\theta^{(i)}$

$$\log l(\theta) = \log l(\theta^{(i)}) + (\theta - \theta^{(i)})^T \nabla \log l(\theta^{(i)}) \quad \text{Taylor expansion}$$

$$+ \frac{1}{2} (\theta - \theta^{(i)})^T [\nabla^2 \log l(\theta^{(i)})] (\theta - \theta^{(i)}) + \text{remainder}$$

$\nabla \log l(\theta^{(i)})$ : gradient at  $\theta^{(i)}$

$\nabla^2 \log l(\theta^{(i)})$ : Hessian matrix at  $\theta^{(i)}$  :  $H_{\theta^{(i)}}$

$$\text{After differentiating: } \nabla \log l(\theta) \approx H_{\theta^{(i)}} (\theta - \theta^{(i)}) + \nabla \log l(\theta^{(i)}) = 0$$

Newton-Raphson method:

$$\text{Initialize } \theta^{(0)} = \hat{\theta}^{(0)}$$

$$1. \theta^{(i+1)} = \theta^{(i)} - [H_{\theta^{(i)}}]^{-1} \nabla \log l(\theta^{(i)}) \quad \text{inverse} \times \text{gradient}$$

until convergence. ( $J$ )

If  $\log l(\theta)$  is quadratic, NR will find solution in one step (i.e. normal distribution).  
In GLM, no are extrema, no local extrema!

The quantity  $E[H]$  (Hessian) determines the sharpness of the peak in the  
likelihood function around its maximum

Alternative to Newton Raphson: METHOD OF SCORING

$$I = E[J|\theta] \quad (\text{Expected Fisher Information})$$

Replaces  $J$  by  $I$  in NR method

20/09/16

## ALSM 2 EXAM NOTES

Binomial

$$P(y|n) = \binom{n}{y} \theta^y (1-\theta)^{n-y} \quad y \in \{0, 1, \dots, n\}$$

Logit ( $\frac{\theta}{1-\theta}$ ) from II - inverse normDistribution:  $= [p+(1-p)]^n = [p]^n = 1$  by binomial expansion

$$\text{Expectation: } \sim \exp[y \log(\theta) + \log(1) + n \log(1-\theta)]$$

$$\text{Momentum: } [y] [y^y (1-\theta)^{n-y}] [y ((1-\theta) \cdot d\ln y)] \quad \theta = y/n$$

$$\text{Expectation Value: } y(\theta) = n \left[ \frac{y}{n} \right] = n \theta \quad n \theta (p+n-q)^n = n \theta$$

## POISSON DISTRIBUTION

$$P(y|n) = \frac{\lambda^y}{y!} e^{-\lambda} \quad \lambda \geq 0 \quad y \in \mathbb{N}$$

Distribution:  $\exp(-\lambda) = \frac{\lambda^0}{0!} \quad (\text{takes exp}) \quad \exp(-\lambda) \exp(\lambda) = 1$ 

$$\text{Expectation: } \exp[y \log(\lambda) - \log(y!) - \lambda]$$

$$\text{Expectation: } \exp(\lambda) = \frac{\lambda^0}{0!} = \lambda^0 e^{\lambda^2/2} e^{-\lambda} = \lambda \exp(-\lambda) \exp(\lambda) = \lambda$$

$$\text{Momentum: } [\lambda^y \frac{d\ln \lambda}{dy}] [y - \lambda] = 0 \quad \lambda = y$$

## Exponential DISTRIBUTION

$$P(y|n) = \theta \exp(-\theta y) \quad y \in \mathbb{R}^+ \quad \theta \in \mathbb{R}^+$$

Distribution:  $\int_0^\infty \theta e^{-\theta y} dy = -[\exp(-\theta y)]_0^\infty = 0 + 1 = 1$ 

$$\text{Expectation: } \exp[y \log(\lambda) - \lambda y]$$

$$\text{Expectation: } [-y \exp(-\lambda y)]_0^\infty + \int_0^\infty \exp(-\lambda y) dy [0 - 0] + \left[ \frac{1}{\lambda} \exp(-\lambda y) \right]_0^\infty = \frac{1}{\lambda}$$

$$\text{Momentum: } \exp(-\lambda y) - \lambda y \exp(-\lambda y) \quad \exp[-\lambda y][1 - \lambda y] \Rightarrow \lambda = y \quad (\text{mean})$$

$$\text{Failure F(y): } \int_y^\infty \lambda \exp(-\lambda z) dz = [-\exp(-\lambda z)]_y^\infty = 1 - \exp(-\lambda y)$$

$$\text{Survival S(y): } 1 - F(y) = \exp(-\lambda y)$$

$$\text{Hazard: } \text{original function / survival fun} \quad \frac{\Delta \exp(-\lambda y)}{\exp(-\lambda y)} = \lambda \quad (\text{memoryless, not dependent on time})$$

## Weibull DISTRIBUTION

$$P(y|\lambda, \theta) = \lambda \theta y^{\lambda-1} \exp[-\theta y^\lambda] \quad y \in \mathbb{R}^+ \quad \theta \in \mathbb{R}^+ \quad \lambda \in \mathbb{R}^+$$

Distribution:  $\int_0^\infty = [-\exp(-\theta y^\lambda)]_0^\infty = 1 + 0 = 1$ 

$$\text{Expectation: ?} \quad \exp(-\theta y^\lambda + \log(\lambda \theta y^\lambda)) = \exp(-\theta y^\lambda + \log(\lambda) + \log(\theta) + (\lambda+1) \log y)$$

$$\text{Expectation: } \left( \frac{y}{\lambda} \right)^{\lambda+1} \Gamma(1 + \frac{1}{\lambda})$$

Momentum: ?

$$\text{Failure: } P(y|y_0) = 1 - \exp[-\theta y_0^\lambda]$$

$$\text{Survival S(y): } 1 - F(y) = \exp[-\theta y^\lambda]$$

$$\text{Hazard: } \frac{P(y|y_0)}{S(y)} = \lambda y^{\lambda-1} \quad (\text{concentrated failure time function of } y)$$

Multi-Nomial Distribution

$$P(y_1, y_2, \dots, y_J, n) = \frac{n!}{y_1! y_2! \dots y_J!} \theta_1^{y_1} \theta_2^{y_2} \dots \theta_J^{y_J}$$

- Generalization of binomial. Binomial when  $J=2$ . When  $J>2$  no. member of experiment family

-  $J$  categories.  $\theta_i$  probability of category,  $\theta_1 + \theta_2 + \dots + \theta_J = 1$

- We can model  $\vec{y}^T$  collected over  $N$  groups into a set of parameters  $\vec{\beta}$

Reference (category) e.g.  $\theta_1$

$$\text{Log } (\theta_j) = [\theta_j / \theta_1] = x^T \beta_j \quad \forall j=2, J \text{ having constraint } \sum_{j=1}^J \theta_j = 1$$

- When estimates  $\hat{\beta}_j$  are computed then  $\hat{\theta}_j = \hat{\theta}_1 \exp(x^T \hat{\beta}_j) \quad \forall j=2, J$

$$\hat{\theta}_1 = \frac{1}{1 + \sum_{j=2}^J \exp(x^T \hat{\beta}_j)}$$

$$\hat{\theta}_j = \frac{\exp(x^T \hat{\beta}_j)}{1 + \sum_{j=2}^J \exp(x^T \hat{\beta}_j)} \quad \forall j=2, J \text{ SOFTMAX FUNCTION} \Rightarrow \text{converting output to probability}$$

below  $x$  (category)  $\rightarrow$  D.F.

Selected model: Number groups  $\times (k_D p(m)-1)$

Additional example

Selected model has  $N(J-1)$  D.F.

20/05/16

## PLSM 2 EXAM NOTES

### GENERAL AND PREDICTIVE ODDS.

$$h_0: \theta \leq 0.5 \quad v \text{ H} \quad \text{if } z_{0.05} \text{ reject } H_0$$

$$P_\theta = S_0^{-\frac{1}{2}} \left( \frac{1}{2} (1/\theta)(1-\theta)^{1/\theta} \right)^{\theta/(1-\theta)} d\theta$$

$$\int_0^{\infty} \left( \frac{1}{2} (1/\theta)(1-\theta)^{1/\theta} \right)^{\theta/(1-\theta)} d\theta$$

### DEVIANCIE

- Null deviance:  $-2 \log \left[ \frac{\text{summed log (observed)}}{\text{null model (predicted)}} \right]$

- Residual Deviance:  $-2 \log \left[ \frac{\text{summed log (residual)}}{\text{sum log (obs)}} \right]$

- Measure of goodness of fit, smaller value better fit

- If model is good deviance follows  $\chi^2$  distribution with  $n-m$  D.F. and 0 D.F.

- If  $D_{\text{model}} < \chi^2_{n-m, \alpha}$  then it is a good model

### AIC

$$AIC = -2 \log h(\beta) + 2p$$

- $p = \dim(\beta) = \text{number of estimated parameters}$

- $\beta$  are the estimated parameters that maximize the likelihood or log likelihood

- $\log h(\beta)$  is maximum value of log likelihood

- Best model is a tradeoff between one that maximizes likelihood with minimum number of parameters

- Select model with lowest AIC

### SURVIVAL ANALYSIS

- Time until failure  $y \in \mathbb{R}^+$

- Probability of failure:  $S(y) = \Pr[Y \geq y] = 1 - F(y)$

- Probability of survival:  $1 - F(y)$

- Hazard function  $\frac{\Pr[y \leq t | y > y]}{\Pr[y = t]}$  chance of failure between time  $y$  and time  $y + \delta y$  given  $y > y$

- $h(y) = \Pr[y|y]/S(y)$

- "Chances of dying in the next minute, given a survival rate and given you have not died before time  $y"$

- Cumulative hazard function  $H(y) = -\log[S(y)]$

- When  $F(y) = S(y) = 0.5$  called the "median point of failure"

- When  $y_i$  is censored, use  $S(y_i)$  instead of  $\Pr[y_i|y_i]$

$$= \prod_{i=1}^N \Pr[y_i|y_i]^{c_i} S(y_i)^{1-c_i}$$

### NEWTON RAPHSON

$$\text{Initial } \theta^{(0)} = \theta^{(0)}$$

$$d^{(0)} < \theta^{(1)} - [H_{\theta^{(0)}}]^{-1} \nabla \log L(\theta^{(0)})$$

Local convergence only one extremum in GRM (no local extrema)

### Method of Scoring

$$I = E[\pi \theta] \quad \text{-Expected Fisher Information}$$

Response  $=[H_{\theta^{(1)}}]$  by  $I$  in NR method

### Limitation of likelihood

-Assumed responses are independent

-Outliers can give a zero all data except one group visibility for parameter  $\beta$

### GLM STUFF ✓

- Bernoulli ✓
  - Binomial ✓
  - Poisson ✓
  - link between Poisson/Binomial ✓
  - Normal/Gaussian
  - Multinomial
  - Survival: Weibull
- Exponential
- Deviance/AIC ✓
  - Binomial Deviance ✓

1. Show this is a sufficient family
2. Show number of exponential family
3. Compute expectation, variance, CDF
4. Compute log likelihood wrt  $\theta$

Basic SLR ✓

Problem Sheet 1 ✓

ANOVA DECOMPOSITION ✓ DF/SS etc

2 ✓

Properties Expectation and Variance ✓

3

Hypothesis testing / Confidence Interval! ✓

4

F-statistic - P value ✓

(I and PI ✓

Matrix Formulation of SLR Model

Multiple Regression

ANOVA

CI PI

Sequential and Partial SS

General linear hypothesis (L matrix)

Orthogonality / Multicollinearity?

