

Jan 26

Table of counts

column factor (response)

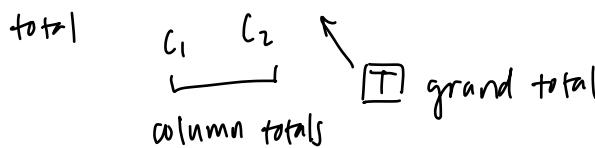
		D C		total	row totals
Row factor (exploratory)	S	n_{11}	n_{12}		
	NS	n_{21}	n_{22}		

S: smoker

NS: nonsmoker

D: disease

C: control



fix row totals first

• Poisson sampling

• Prospective product binomial

• Multinomial sampling

• Randomized experiment

• Retrospective product binomial

)

fix column totals first

• Hypothesis of homogeneity : $H_0: \pi_1 = \pi_2, \frac{w_1}{w_2} = 1$

• Hypothesis of independence : $H_0: \pi_1 = \pi_2,$

odds

odds ratio

$\pi_1: P(\text{observe D in first row})$

$$w_1 = \frac{\pi_1}{1-\pi_1}$$

$$\frac{w_1}{w_2}$$

$\pi_2: P(\text{observe D in second row})$

$$w_2 = \frac{\pi_2}{1-\pi_2}$$

* can be done when $\frac{w_1}{w_2} = 1$

b/c otherwise you don't know true proportions of D vs. C in population

$$\frac{w_1}{w_2} = \frac{\frac{\pi_1}{1-\pi_1}}{\frac{\pi_2}{1-\pi_2}} = \frac{\pi_1}{\pi_2} \cdot \frac{1-\pi_2}{1-\pi_1}$$

$$= \frac{\frac{\pi_1}{\pi_2}}{\frac{1 - \pi_1}{1 - \pi_2}} >$$

both of these ratios can be correctly estimated using retrospective product binomial

$$H_0: \pi_1 = \pi_2$$

test proportions:

z-test:

$$\hat{\pi}_1 = \frac{n_{11}}{R_1}, \quad \hat{\pi}_2 = \frac{n_{21}}{R_2}$$

$$\hat{\pi}_c = \frac{n_{11} + n_{21}}{R_1 + R_2} \quad (\text{assume two rows have common } \pi)$$

$$E(\hat{\pi}_1 - \hat{\pi}_2) = \pi_c - \pi_c, \quad \text{Var}(\hat{\pi}_1 - \hat{\pi}_2) = \frac{\pi_c(1-\pi_c)}{R_1} + \frac{\pi_c(1-\pi_c)}{R_2}$$

$$z = \frac{\hat{\pi}_1 - \hat{\pi}_2}{\sqrt{\frac{\hat{\pi}_c(1-\hat{\pi}_c)}{R_1} + \frac{\hat{\pi}_c(1-\hat{\pi}_c)}{R_2}}} \sim N(0, 1)$$

chi-squared test:

$$\chi^2 = \sum_{ij} \frac{(Y_{obs} - Y_{exp})^2}{Y_{exp}}$$

$$N \chi^2_{df} = (p-1)(q-1)$$

p rows, q columns

		Obs		Exp	
		n_{11}	n_{12}	R_1	$\frac{R_1 C_1}{T}$
		n_{21}	n_{22}	R_2	$\frac{R_2 C_2}{T}$
C ₁	C ₂			T	

Fisher's Exact test:

z-test and χ^2 rely on asymptotic distributions. Fisher's does not rely on asymptotic assumptions (can be used for small sample size) (hypergeometric dist.)

$$\Pr(n_{11} = k) = \frac{\binom{R_1}{k} \binom{R_2}{C_1 - k}}{\binom{T}{C_1}}$$

can use this to compute a one-tailed p-value

If $n_{11} > \frac{R_1 G_1}{T}$ then $p\text{-val} = \sum_{k=n_{11}}^{\min(R_1, G_1)} \Pr(n_{11} = k)$

If $n_{11} < \frac{R_1 G_1}{T}$ then $p\text{-val} = \sum_{k=0}^{n_{11}} \Pr(n_{11} = k)$

Binary outcome

$\text{Ber}(\pi)$ Binomial data:

$$Y \sim \text{Bin}(m, \pi)$$

/ binomial denominator

$$(i) Y = Y_1 + \dots + Y_n$$

$$(ii) \begin{aligned} Y_1 &\sim \text{Pois}(m_1) & \Pr(Y_1 = y | Y_1 + Y_2 = m) = \binom{m}{y} \pi^y (1-\pi)^{m-y} \\ Y_2 &\sim \text{Pois}(m_2) & \pi = \frac{m_1}{m_1 + m_2} \end{aligned}$$

Moments and cumulants

For Bernoulli:

$$\begin{aligned} \text{MGF: } M_Y(\gamma) &= E e^{\gamma Y} \\ &= (1-\pi) + \pi e^\gamma \end{aligned}$$

Cumulant generating function:

$$\begin{aligned} K_Y(\gamma) &= \log M_Y(\gamma) \\ &= \log((1-\pi) + \pi e^\gamma) \end{aligned}$$

1st cumulant = mean

2nd cumulant = 2nd central moment

3rd cumulant = 3rd central moment

No corresponding relationship for 4th and beyond.

For Binomial:

$$\text{MGF : } M_Y(\gamma) = (1-\pi + \pi e^\gamma)^m$$

Momulant generating function:

$$K_Y(\gamma) = m \log(1 - \pi + \pi e^\gamma)$$

Taylor expansion of CGF:

$$K_Y(\gamma) = K_Y(0) + K_Y'(0)\gamma + \frac{1}{2!} K_Y''(0)\gamma^2 + \dots$$

$$K_Y'(0) = \left. \frac{m\pi e^\gamma}{1 - \pi + \pi e^\gamma} \right|_{\gamma=0} = m\pi$$

$$\begin{aligned} K_Y''(0) &= m\pi \left(\frac{e^\gamma}{1 - \pi + \pi e^\gamma} - \frac{\pi e^{2\gamma}}{(1 - \pi + \pi e^\gamma)^2} \right) \\ &= m\pi \left. \frac{(1 - \pi)e^\gamma}{(1 - \pi + \pi e^\gamma)^2} \right|_{\gamma=0} = m\pi(1 - \pi) \end{aligned}$$

$$K_Y'''(0) = m\pi(1 - \pi)(1 - 2\pi)$$

$$K_Y^{(4)}(0) = m\pi(1 - \pi)[1 - 6\pi(1 - \pi)] \quad \begin{matrix} \rightarrow \\ \text{derive as} \\ \text{exercise} \end{matrix}$$

Jan 31

Cumulant generating function:

$$K_Y(t) = \log M_Y(t)$$

$$\text{Binomial: } K_Y(t) = m \underbrace{\log \{1 - \pi + \pi e^t\}}_{\text{CGF for Bernoulli}}$$

Normal limit

$$Z = \frac{Y - m\pi}{\sqrt{m\pi(1-\pi)}} \xrightarrow{m \rightarrow \infty} N(0, 1)$$

Standard Normal CGF: $\sqrt{t^2/2}$ (derive as exercise)

$$0 \ 1 \ 0 \ \dots \ 0$$

scale y :

$$\begin{aligned} \text{original CGF: } & E e^{ty} \\ &= K_0 + K_1 t + \frac{K_2}{2!} t^2 + \dots \end{aligned}$$

after scaling: $\log E e^{tcy}$, $\tilde{y} = cy$

$$= K_0 + K_1 ct + \frac{K_2}{2!} (ct)^2 + \dots \quad \sqrt{m} \text{ serves as scaling factor}$$

$$K_0 \quad cK_1 \quad c^2 K_2 \quad \dots$$

$$N(0, 1): \quad \begin{matrix} K_1 & K_2 & K_3 \\ 0 & 1 & 0 \dots 0 \end{matrix}$$

$$Z: \quad \begin{matrix} 0 & 1 & 0 \left(\frac{1}{\sqrt{m}}\right) & 0 \left(\frac{1}{m}\right) & 0 \left(\frac{1}{m^{5/2}}\right) \\ \uparrow & \downarrow & \downarrow & \downarrow & \downarrow \end{matrix}$$

slowest converging term

convergence rate = \sqrt{m}

Poisson limit:

Poisson CGF: $m(e^t - 1)$

all cumulants determined by m (show w/ Taylor expansion)

$$Y: m \log(1 - \pi + \pi e^t)$$

consider the following case:

$$m \rightarrow \infty, \pi \rightarrow 0 \text{ s.t. } m\pi \rightarrow \mu$$

$$m \log(1 - \pi + \pi e^t)$$

$$\rightarrow \frac{m}{\pi} \log(1 + \pi(e^t - 1))$$

$$= \frac{m}{\pi} \left(\pi(e^t - 1) - \frac{1}{2} \pi^2 (e^t - 1)^2 + \frac{1}{3} \pi^3 (e^t - 1)^3 \dots \right)$$

$$= m(e^t - 1) - \frac{1}{2} m \pi (e^t - 1)^2 + \frac{1}{3} m \pi^2 (e^t - 1)^3 \dots$$

$\overbrace{\quad}$
Poisson CGF

$$\underbrace{\frac{m}{m}}_{O(\frac{1}{m})} \downarrow 0$$

$$\underbrace{O\left(\frac{1}{m^2}\right)}_{0} \downarrow 0$$

\Rightarrow Binomial can be approximated w/ Poisson

convergence rate = $1/m$ (faster than \sqrt{m})

Models for Binary / Binomial Responses

$$y \quad \text{Ber}(\pi) \quad 0 \leq \pi \leq 1$$

$\text{Binomial}(m, \pi)$

- need to choose a link function for $X\beta$, π .
can't use identity b/c $X\beta$ is not bounded.

(1) logit / logistic function

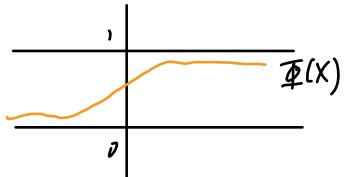
$$g_1(\pi) = \log \frac{\pi}{1-\pi}$$

inverse CDF for
logistic distribution:

$$\frac{e^\pi}{1-e^\pi}$$

(2) Probit / inverse Normal function

$$g_2(\pi) = \Phi^{-1}(\pi)$$



(3) complementary log-log function

$$g_3(\pi) = \log \{-\log(1-\pi)\}$$

extreme value

distribution CDF:

$$1 - e^{-e^\pi}$$

(4) log-log link

$$g_4(\pi) = -\log \{-\log(\pi)\}$$

extreme value CDF:

$$e^{-e^\pi}$$

· g_1 and g_2 are symmetric around $\frac{1}{2}$.

$$g_1(\pi) = g_1(1-\pi)$$

$$g_2(\pi) = g_2(1-\pi)$$

· g_3 and g_4 are related.

$$g_3(\pi) = -g_4(1-\pi)$$

· g_1 is most commonly used.

Logistic Regression

$$\log \frac{\pi}{1-\pi} = X\beta \quad Y \sim \text{Bin}(m, \pi)$$

$$\text{ex: } X\beta = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

1. parameter interpretation:

β_0 : log odds of getting disease when $X_1, X_2 = 0$

$$\frac{e^{\beta_0}}{1 + e^{\beta_0}} = \pi$$

β_1 : hold all other variables constant, change in log odds w/ one unit increase of X_1

$$\log \frac{\pi_1}{1-\pi_1} = \beta_0 + \beta_1(X_1+1) + \beta_2 X_2$$

$$\log \frac{\pi_1}{1-\pi_1} = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$$\log \frac{\pi_1}{1-\pi_1} - \log \frac{\pi_2}{1-\pi_2} = \beta_1$$

$$= \log \frac{\frac{\pi_1}{1-\pi_1}}{\frac{\pi_2}{1-\pi_2}}$$

β_1 characterizes a log-odds ratio

$$e^{\beta_1} \rightarrow \text{odds ratio}$$

- Randomness incorporated into distributional assumption, no error term.

Retrospective trial:

β_1 and β_2 still reflect population behavior.

Assume this true model for the population:

$$\text{logit } P(D|X) = \alpha + \beta^T X \quad \bar{D} = \text{no disease}$$

D = disease

but data collected is not representative of pop.

$Z=1$ patient sampled

$Z=0$ patient not sampled

$$\pi_1 = P(Z=1|D) \quad \pi_0 = P(Z=1|\bar{D})$$

fit: logit $P(D|X, Z=1)$

$$P(D|X, Z=1) = \frac{P(D, Z=1|X)}{P(Z=1|X)}$$

$$= \frac{P(Z=1|D, X) P(D|X)}{P(Z=1|D, X) P(D|X) + P(Z=1|\bar{D}, X) P(\bar{D}|X)}$$

$$= \frac{\pi_1 P(D|X)}{\pi_1 P(D|X) + \pi_0 P(\bar{D}|X)}$$

$$= \pi_1 \frac{e^{\alpha + \beta^T X}}{1 + e^{\alpha + \beta^T X}}$$

$$\frac{\pi_1 \frac{e^{\alpha + \beta^T X}}{1 + e^{\alpha + \beta^T X}} + \pi_0 \frac{1}{1 + e^{\alpha + \beta^T X}}}{\pi_1 \frac{e^{\alpha + \beta^T X}}{1 + e^{\alpha + \beta^T X}} + \pi_0}$$

$$= \frac{\pi_1 e^{\alpha + \beta^T X}}{\pi_1 e^{\alpha + \beta^T X} + \pi_0}$$

$$= \frac{\frac{\pi_1}{\pi_0} e^{\alpha + \beta^T X}}{\pi_1 e^{\alpha + \beta^T X} + 1}$$

$$P(D|X, \mathcal{Z}=1) = \frac{e^{\log \frac{\pi_1}{\pi_0} + \alpha + \beta^T X}}{e^{\log \frac{\pi_1}{\pi_0} + \alpha + \beta^T X} + 1}$$

$$\text{logit } P(D|X, \mathcal{Z}=1) = \underbrace{\log \frac{\pi_1}{\pi_0} + \alpha}_{\text{offset}} + \beta^T X$$

without sampling probabilities you can only talk about
 β 's (can't interpret intercept)

$$\alpha^* = \alpha + \log \frac{\pi_1}{\pi_0}$$

log-likelihood:

$$l(\pi, y) = \sum_{i=1}^n \log \left(\frac{m_i}{y_i} \right) \pi_i^{y_i} (1 - \pi_i)^{m_i - y_i}$$

$$= \sum_{i=1}^n \left[y_i \log \frac{\pi_i}{1 - \pi_i} + m_i \log (1 - \pi_i) \right] + \text{const}$$

logit of x_{ij}
anything w/o
unknown
param

$$l(\beta, y) = \sum_{i=1}^n \sum_j y_{ij} x_{ij} \beta_j - \sum_{i=1}^n m_i \log \left(1 + e^{\sum_j x_{ij} \beta_j} \right)$$

$$= \sum_{i=1}^n y_{ij} x_{i\cdot}^\top \beta - \sum_{i=1}^n m_i \log \left(1 + e^{x_{i\cdot}^\top \beta} \right)$$

$$= y^\top X \beta - \sum_{i=1}^n m_i \log \left(1 + e^{x_{i\cdot}^\top \beta} \right)$$

gulf. stat:
 $x_{i\cdot}^\top y$

$$\text{MLE: } l'(\beta, y) = 0$$

$$x^T y - \sum_i m_i \underbrace{\frac{x_i}{1 + e^{x_i^T \beta}}}_{\text{not linear in terms of } \beta} = 0$$

not linear in terms of β

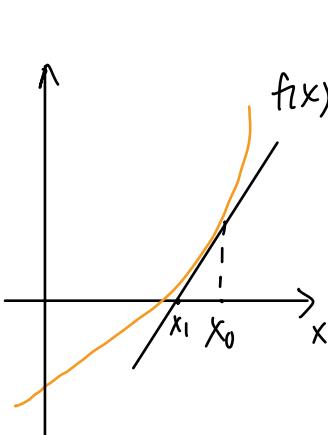
Feb 2

Hw: Read chp. 5 of text

logistic regression MLE

$$\begin{aligned} l(\beta, y) &= \sum_{i=1}^n \left[y_i \log \frac{\pi_i}{1-\pi_i} + m_i \log (1-\pi_i) + c \right] \\ &= y^T X \beta - \sum_{i=1}^n m_i \log (1 + e^{x_i^T \beta}) \end{aligned}$$

$\frac{\partial l}{\partial \beta} = 0$, solve using Newton-Raphson



$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

$$x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$$

:

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

converges to root of f(x)

after iteratively running update

$$y - f(x_0) = f'(x_0)(x - x_0)$$

$$y = f(x_0) + f'(x_0)(x - x_0)$$

$$x = x_0 - \frac{f(x_0)}{f'(x_0)}$$

if β is multi-dimensional:

$$\beta_1 = \beta_0 - \left[\underbrace{\mathbb{E} l''(\beta_0)}_{\text{Fisher info}} \right]^{-1} \underbrace{l'(\beta_0)}_{\text{score}}$$

↑
usually take expectation b/c
 X is an RV (Fisher scoring)

$$\vdots$$

$$\beta_{t+1} = \beta_t - \left[\mathbb{E} l''(\beta_t) \right]^{-1} l'(\beta_t)$$

Deriving a concrete form:

$$(1) \quad \frac{\partial l}{\partial \beta} = X^T y - \sum_{i=1}^n \frac{m_i e^{x_i^T \beta}}{1 + e^{x_i^T \beta}} x_i$$

$$= X^T y - (X_1 \ X_2 \dots \ X_n) \begin{pmatrix} m_1 \pi_1 \\ m_2 \pi_2 \\ \vdots \\ m_n \pi_n \end{pmatrix}$$

$$= X^T (y - M) \quad X^T y \text{ is a sufficient statistic for } \beta, \mathbb{E}(X^T y) = X^T M$$

Another way:

$$\frac{\partial l}{\partial \beta} = \begin{bmatrix} \frac{\partial l}{\partial \beta_1} \\ \vdots \\ \frac{\partial l}{\partial \beta_p} \end{bmatrix}$$

$$\frac{\partial \ell}{\partial \beta_r} = \sum_{i=1}^n \frac{\partial \ell}{\partial \pi_i} \frac{\partial \pi_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_r}$$

\$\pi_i\$ is a function
of linear predictor
(link function)
(\$\eta_i\$ is linear predictor)

$$= \sum_{i=1}^n \frac{y_i - m_i \pi_i}{\pi_i(1-\pi_i)} \pi_i(1-\pi_i) x_{ir}$$

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

$$= \sum_{i=1}^n (y_i - m_i \pi_i) x_{ir} \quad \frac{d\pi_i}{d\eta_i} = \frac{e^{\eta_i}}{1+e^{\eta_i}} - \frac{e^{2\eta_i}}{(1+e^{\eta_i})^2}$$

$$= \pi_i(1-\pi_i)$$

r^{th} entry of the gradient vector

Second derivative: (only π_i involves β in $\frac{\partial \ell}{\partial \beta_r}$)

$$- \mathbb{E} \frac{\partial^2 \ell}{\partial \beta_r \partial \beta_s} = \mathbb{E} - \sum_{i=1}^n -m_i x_{ir} \frac{d\pi_i}{d\eta_i} \frac{\partial \eta_i}{\partial \beta_s}$$

same result = $\mathbb{E} \sum_{i=1}^n m_i x_{ir} \pi_i(1-\pi_i) x_{is}$

with or

without taking expectation = $\sum_{i=1}^n m_i \pi_i(1-\pi_i) x_{ir} x_{is}$

$$-\mathbb{E} \frac{\partial^2 \ell}{\partial \beta^2} = X^T W X$$

where $W = \text{diag} \left\{ \underbrace{m_i \pi_i (1 - \pi_i)}_{\text{variance for the } i^{\text{th}} \text{ subject}} \right\}$

So iterations for Newton-Raphson look like:

$$\beta_{t+1} = \beta_t + (X^T W X)^{-1} X^T (y - m)$$

↑
function of β_t

Another perspective: iterated least squares

$$\begin{aligned}\beta_1 &= \beta_0 + (X^T W X)^{-1} X^T (y - m) \\ &= \beta_0 + (X^T W X)^{-1} X^T W W^{-1} (y - m) \\ &= (X^T W X)^{-1} \left\{ (X^T W X) \beta_0 + X^T W W^{-1} (y - m) \right\} \\ &= (X^T W X)^{-1} X^T W (X \beta_0 + W^{-1} (y - m))\end{aligned}$$

$$\tilde{z} = \underbrace{\mathbf{x}\beta_0}_\text{linear predictor} + \underbrace{W^{-1}(y - \mu)}_\text{observed - expected error term} \rightarrow \begin{array}{l} \text{looks like response variable} \\ \text{in linear models} \end{array}$$

$$\beta_1 = (\mathbf{x}^T W \mathbf{x})^{-1} \mathbf{x}^T W \tilde{z}$$

weighted least squares estimate for β

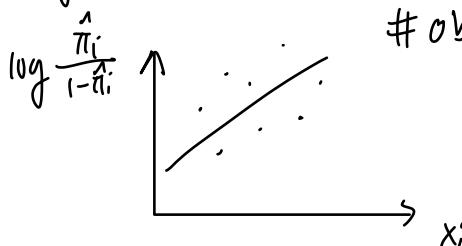
$$\tilde{z} = \mathbf{x}\beta + \varepsilon, \quad \varepsilon \sim N(0, \propto W^{-1})$$

alternatively: $W^{1/2} \tilde{z} = W^{1/2} \mathbf{x}\beta + W^{1/2} \varepsilon$

$$\tilde{z} = \tilde{\mathbf{x}}\beta + \tilde{\varepsilon} \sim N(0, \propto I)$$

$$\begin{aligned} (\tilde{\mathbf{x}}\tilde{\mathbf{x}})^{-1} \tilde{\mathbf{x}}\tilde{\varepsilon} &= (\mathbf{x}^T W^{1/2} W^{1/2} \mathbf{x})^{-1} \mathbf{x}^T W^{1/2} W^{1/2} \varepsilon \\ &= (\mathbf{x}^T W \mathbf{x})^{-1} \mathbf{x}^T W \tilde{z} \end{aligned}$$

estimating β_0 : have to assume that



obs > # params

Another perspective: link function

higher order terms

$$g(y) = g(\mu) + g'(\mu)(y-\mu) + O(\epsilon)$$

$$= X\beta + \frac{d\eta}{d\mu}(y-\mu)$$

based on model assumptions

So adjusted response \tilde{z} is using Taylor expansion to approximate link function applied to observed data.

Also called iteratively re-weighted least squares (IRWLS)

Inference :

- for MLEs, $\text{Var } \hat{\beta}_{\text{mle}} \rightarrow$ inverse of asymptotically fisher info
 - When the algorithm converges, we already obtain the fisher info.

$$\text{Var}(\hat{\beta}) \cong (X^T W X)^{-1}$$

$$\beta_j \text{ logit}(\pi) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$z_j = \frac{\hat{\beta}_j - \beta_{j0}}{\sqrt{\text{var}(\hat{\beta}_j)}} \quad \text{Wald test}$$

/

asymptotic variance, var_{jj}

Deviance

$$\text{Residual deviance} \triangleq 2(\ell_{\max} - \ell)$$

maximum possible ℓ w/
fully saturated model

(each observation has its own π_i
w/o being constrained by shared
function for all obs)

Binomial logistic

$$D(y, \hat{\pi}) = 2 \left\{ \ell(\hat{\pi}, y) - \ell(\pi, y) \right\}$$

$$= 2 \sum_i \left\{ y_i \log \frac{y_i}{m_i \hat{\pi}_i} + (m_i - y_i) \log \left(\frac{1 - \frac{y_i}{m_i}}{1 - \hat{\pi}_i} \right) \right\}$$

small deviance \rightarrow model sufficient to describe the data.

$$D \sim \chi^2_{n-p} \text{ (under certain conditions)}$$

can use χ^2 to determine if D is big or small

$$\text{mean } \chi^2_{n-p} = n-p$$

cannot be used for binary outcome

Drop-in-deviance test:

for comparing a null model w/ an alternative model.

H_0 : simpler model

H_1 : more complex model

Drop-in-deviance statistic:

$$D(y_i; \hat{\mu}_0) - D(y_i; \hat{\mu}_1)$$

↑
MLE from null model

MLE from alternative model

likelihood ratio

$$= 2(\ell_{\max} - \ell_0) - 2(\ell_{\max} - \ell_1) = 2(\ell_1 - \ell_0)$$

from likelihood ratio test \rightarrow $\sim \chi^2_{p_1 - p_0}$

• this is an asymptotic result.

$m \rightarrow \infty$ p_1, p_0 fixed.

either
are
OK

as $n \rightarrow \infty$, per-parameter information increases
(different from deviance test)

Feb 9

- models for polytomous data
- polytomous response: has several categories

ex: Blood types (unstructured) $\leftarrow \pi_A \pi_B \pi_{AB} \pi_O$
levels of education (ordered)

$$\{ \underbrace{\pi_1, \dots, \pi_K}_{\text{response probabilities (not ordered)}}, \sum_{i=1}^K \pi_i = 1$$

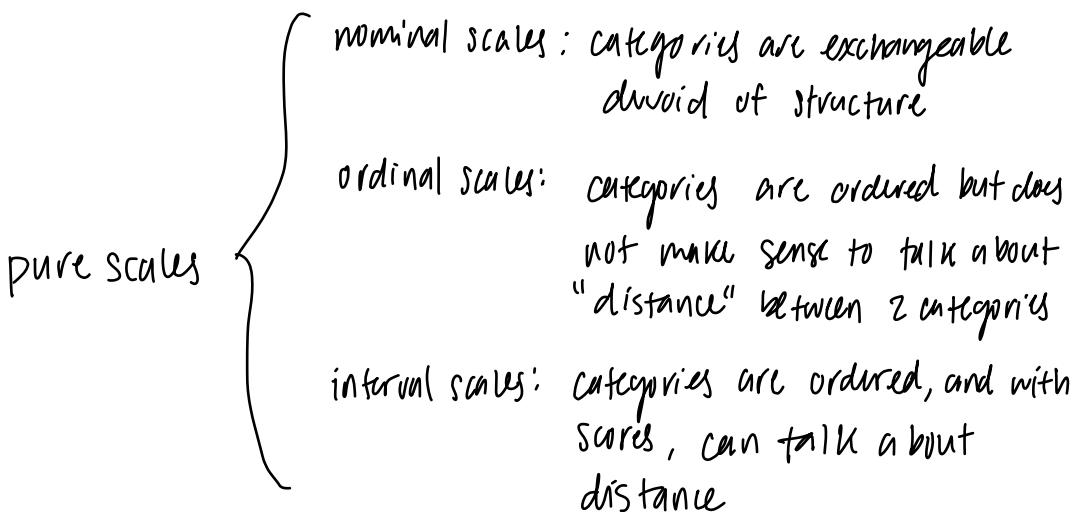
$$\gamma_1 = \pi_1, \gamma_2 = \pi_1 + \pi_2, \gamma_3 = \pi_1 + \pi_2 + \pi_3 \dots$$

for ordered responses

$$\text{(i.e. } \gamma_2 = P(\text{middle + high school}) \text{)}$$

} cumulative response probabilities

Measurement scales



Compound scales:

ex: bivariate response $\begin{pmatrix} y \\ z \end{pmatrix} \rightarrow$ ordinal
 \rightarrow binary

binary response $\left\{ \begin{array}{l} \text{nominal} \\ \text{ordinal} \end{array} \right.$

Modelling multi-category responses

$$y_i \sim X_i$$

$$\begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix} \quad / \quad \# \text{ trials}$$

Nominal. $y \sim \text{multinomial}(m, \pi = (\pi_1, \dots, \pi_k))$

$$\pi \sim X \text{ and } \sum_{i=1}^k \pi_i = 1$$

Commonly used model: baseline category logit model
(BCLM)

$$\log \frac{\pi_j^{(x)}}{\pi_k^{(x)}} = \theta_j + \beta_j^T x \quad j = 1, \dots, k-1$$

()
intercept vector containing regression coeffs

similar to binomial logistic regression

will have $K-1$ models. This set of models defines the BCLM.

interpretation:

$$\beta_j = \begin{pmatrix} \beta_{j1} \\ \vdots \\ \beta_{jp} \end{pmatrix} \rightarrow \beta_{jq}$$

$$x \rightarrow x + \begin{pmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{pmatrix} \rightarrow \beta_{jq} = x + e_q \triangleq x'$$

$$\log \frac{\pi_j(x)}{\pi_K(x)} = \theta_j + \beta_j^T x$$

$$\log \frac{\pi_j(x')}{\pi_K(x')} = \theta_j + \beta_j^T x + \beta_j^T e_q$$

$$\log \frac{\pi_j(x')}{\pi_K(x')} - \log \frac{\pi_j(x)}{\pi_K(x)} = \beta_j^T e_q = \beta_{jq}$$

difference between
log odds ratio
between two
categories

w/ 1-unit increase
in q^{th} covariate

model for ordinal scales

Cumulative logit model

$$\text{logit } (\gamma_j) = \theta_j - \beta^T x$$

$$\log \frac{\gamma_j}{1-\gamma_K} = \log \frac{P(Y \leq j)}{P(Y > j)} = \log \frac{\pi_1 + \dots + \pi_j}{\pi_{j+1} + \dots + \pi_K}, \quad j=1, \dots, K-1$$

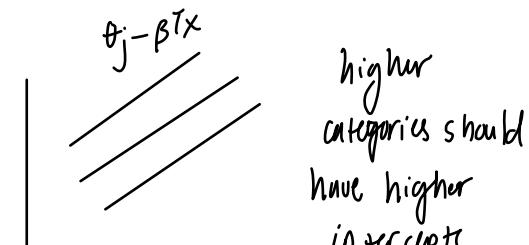
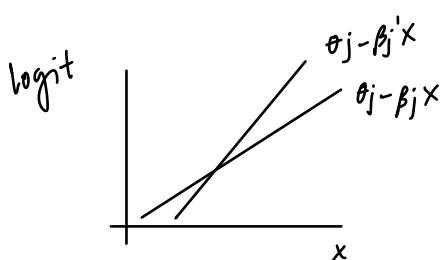
$$\gamma_K = 1$$

have an intercept and vector of coeffs for each model

this and BCLM have the same complexity + same response probabilities after fitting

This model is naive b/c it doesn't consider ordering.

$$\gamma_1 \leq \gamma_2 \leq \dots \leq \gamma_{K-1} \leq \gamma_K$$



same slope for all categories

$$\text{logit } \gamma_j = \theta_j - \beta^T x$$

$$\theta_1 \leq \theta_2 \leq \dots \leq \theta_{K-1}$$

} modified cumulative logit model
(proportional odds model)

$$\log \frac{\gamma_j(x_1)}{1 - \gamma_j(x_1)} \quad \theta_j = \beta^T x_1$$

$$\log \frac{\gamma_j(x_2)}{1 - \gamma_j(x_2)} \quad \theta_j = \beta^T x_2$$

difference:

$$\frac{\frac{\gamma_j(x_1)}{1 - \gamma_j(x_1)}}{\frac{\gamma_j(x_2)}{1 - \gamma_j(x_2)}} = e^{-\beta^T(x_1 - x_2)}$$

education level example:

$$\frac{\frac{M}{H}}{1 - \frac{M}{H}} = \frac{\frac{C}{G}}{1 - \frac{C}{G}}$$

$$\frac{\frac{\gamma_M}{1 - \gamma_M}}{1 - \frac{\gamma_M}{1 - \gamma_M}} = \frac{\frac{\gamma_C(x_1)}{1 - \gamma_C(x_1)}}{1 - \frac{\gamma_C(x_1)}{1 - \gamma_C(x_1)}}$$

different categories show the same behavior in the odds ratio as a function of the covariates

ex: integrating surveys

y: (1) liberal moderate conservative

(2)

very liberal liberal moderate conservative very conservative

↑
 β

↑
 β

can use same proportional odds model to describe both datasets and then integrate them.

proportional odds model can also be viewed from a latent variable perspective.

Z: unobserved RV. Assume $Z - \beta^T x \sim$ standard logistic distribution

$$CDF = \frac{1}{1+e^{-x}} = \frac{e^x}{1+e^x}$$

$$\theta_1 \leq \theta_2 \leq \dots \leq \theta_{j-1}$$

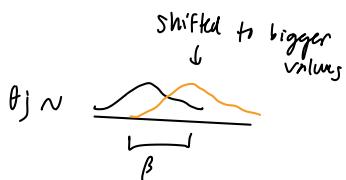
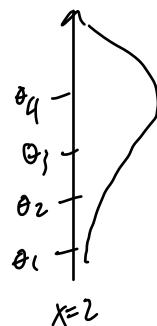
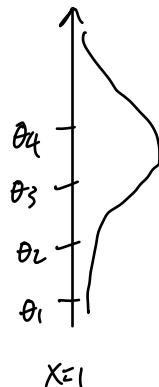
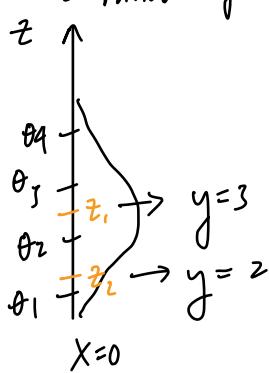
(i) Generate Z

(ii) If $\theta_{j-1} \leq z \leq \theta_j$, then $y = j$

$$\begin{aligned} P(Y \leq j) &= P(Z \leq \theta_j) = P(Z - \beta^T x \leq \theta_j - \beta^T x) \\ &= \frac{e^{\theta_j - \beta^T x}}{1 + e^{\theta_j - \beta^T x}} \end{aligned}$$

$$\Rightarrow \text{logit } P(Y \leq j) = \theta_j - \beta^T x$$

sample z randomly



x increases \rightarrow probability of seeing data from higher categories increases

If the distribution becomes complementary log-log:

extreme \leftrightarrow comp. log-log

(i) generate $z \sim \text{extreme dist}$

$$\log(-\log(1 - r_j)) = \theta_j - \beta^T x \quad] \text{ proportional hazards model}$$

normal \leftrightarrow probit

Feb 14

General Linear Models

- consists of 3 core components:

1. Random component

2. Systematic component

$$\text{linear predictor } \eta = X\beta = \sum_{j=1}^p x_j \beta_j$$

3. Link $g(\mu) = \eta$ usually mean of y or something that describes the location of the response

Linear regression: $\epsilon \sim N(0, \sigma^2 I)$, $y = X\beta + \epsilon \sim N(X\beta, \sigma^2 I)$

No error term in logistic regression b/c you're looking at a parameter

- Exponential family (EF):

$$f(x) = h(x) \exp\{\eta(\theta) T(x) - A(\theta)\}$$

special case: Natural EF (NEF):

$$f(x) = h(x) \exp\{\theta x - A(\theta)\}$$

$$\text{NEF} \subset \text{EF}$$

For EF, If $n(\theta) = \theta$, then this exponential family is in canonical form.

Ex: Normal distribution

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{-x^2+2x\mu+\mu^2}{2\sigma^2}}$$

(1) when σ^2 is known:

$$h(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}} \quad n(\mu) = \frac{\mu}{\sigma}$$

$$T(x) = \underbrace{x/\sigma}_{\text{suff. stat}} \quad A(\mu) = \frac{\mu^2}{2\sigma^2}$$

(2) When σ^2 is unknown:

$$h(x) = \frac{1}{\sqrt{2\pi}} \quad n = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2} \right)$$

$$T(x) = (x, x^2)^T \quad A(\eta) = \frac{\mu^2}{2\sigma^2} + \log|\sigma|$$

- density can always be written in canonical form
- canonical form is not unique.

Ex: Beta distribution

$$\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

$$= \exp \left\{ (\alpha-1) \log x + (\beta-1) \log(1-x) + \log \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \right\}$$

- EF but not NEF.

Ex: log normal distribution

$$\frac{1}{\sqrt{2\pi}\sigma x} e^{-\frac{(\log x - \mu)^2}{2\sigma^2}}$$

- EF but not NEF

$$f_Y(y; \theta, \phi) = \exp \left\{ \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right\}$$

$$\begin{array}{ccc} a(\phi) & b(\theta) & c(y, \phi) \\ & \underbrace{}_{\text{cumulant function}} & \end{array}$$

For now, ϕ is assumed to be known
 θ is called canonical parameter.

Log-likelihood:

$$\begin{aligned} \ell(\theta, \phi; y) &= \log f_Y(y; \theta, \phi) \\ &= \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \end{aligned}$$

Note:

$$E \frac{\partial \ell}{\partial \theta} = 0, \quad E \frac{\partial^2 \ell}{\partial \theta^2} + E \left(\frac{\partial \ell}{\partial \theta} \right)^2 = 0$$

$$\Rightarrow E \frac{\partial \ell}{\partial \theta} = E \frac{y - b'(\theta)}{a(\phi)} = 0$$

$$\Rightarrow \underbrace{E y}_{\sim} = b'(\theta)$$

$$\Rightarrow E \frac{\partial \ell^2}{\partial \theta^2} + E \left(\frac{\partial \ell}{\partial \theta} \right)^2 = E \frac{-b''(0)}{a(\phi)} + E \left(\frac{y - b'(\theta)}{a(\phi)} \right)^2$$

$$\underbrace{\frac{d \mu}{d \theta}}_{\sim} = \frac{-b''(\theta)}{a(\phi)} + \frac{\text{var}(y)}{[a(\phi)]^2} = 0$$

$$\Rightarrow \text{var}(y) = a(\phi) b''(\theta)$$

$$= a(\phi) \underbrace{V(m)}_{\text{variance function.}}$$

ϕ : dispersion parameter

$a(\phi)$: a priori chosen function, commonly $a(\phi) = \frac{\phi}{w}$

w : prior weight

Example:

① Normal

$$f_y(y; \theta, \phi) = \frac{1}{\sqrt{2\pi\phi^2}} e^{-\frac{(y-\mu)^2}{2\phi^2}}$$

$$= \exp \left\{ -\frac{y^2 - 2\mu y + \mu^2}{2\phi^2} - \frac{1}{2} \left\{ \frac{y^2}{\phi^2} + \log 2\pi\phi^2 \right\} \right\}$$

$$\theta = \mu \quad b(\theta) = \theta^2/2$$

$$\phi = \sigma^2 \quad a(\phi) = \phi$$

$$c(y, \theta) = -\frac{1}{2} \left(\frac{y^2}{\theta^2} + \log 2\pi\theta^2 \right)$$

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

$$b''(\theta) = 1 = V(m)$$

$$\text{var}(Y) = a(\phi) V(\mu) = \sigma^2$$

② Poisson

$$f_Y(y; \theta, \phi) = \frac{\mu^y e^{-\mu}}{y!}$$

$$= \exp\{y \log \mu - \mu - \log y!\}$$

$$\begin{aligned} \theta &= \log \mu & \phi &= 1 \\ b(\theta) &= \mu = e^\theta & a(\phi) &= 1 \end{aligned}$$

$$b'(\theta) = e^\theta = \mu = E(Y)$$

$$b''(\theta) = e^\theta = \mu = V(\mu) \Rightarrow \text{var}(y) = a(\phi) b''(\theta) = \mu = E(Y)$$

Link function

common links:

complementary link

① Identity e.g. Normal



② Binomial : logit = $\log\left(\frac{\pi}{1-\pi}\right) = \eta$

$$\text{probit} = \eta = \Phi^{-1}(\pi)$$

$$\text{complementary log-log} = \eta = \log\{-\log(1-\pi)\}$$

$$\textcircled{3} \text{ Poisson: } \log n = \log m$$

- canonical link: related to the canonical parameter. transforms location param into canonical param.
- using canonical link means $X^T y$ is a sufficient statistic.

$$f_Y(y; \theta, \phi) = \frac{m^y e^{-m}}{y!}$$
$$= \exp\{y \log m - m - \log y!\}$$

$$\log m = X\beta$$

$y\theta = yX\beta$ if we use canonical link

$y h(X\beta)$ otherwise: $y(\theta) = y g^{-1}(X\beta)$

Common link functions:

① Normal:

$$\eta = \mu$$

② Poisson:

$$\eta = \log \mu$$

③ Binomial:

$$\eta = \log \frac{\pi}{1-\pi}$$

④ Gamma:

$$\eta = \frac{1}{\mu}$$

⑤ Inverse Gaussian:

$$\eta = \frac{1}{\mu^2}$$

Algorithm for fitting GLM:

(1) iterated reweighted least squares perspective
(IRLS)

Basic idea: Fit a regression using a linearized form of link function $g(y)$

$$g(y) = g(\mu) + g'(\mu)(y - \mu)$$

$$= g(\mu) + \frac{d\eta}{d\mu}(y - \mu)$$

linear approx.

w/ Taylor expansion

$$z \stackrel{\Delta}{=} X\beta + \frac{d\eta}{dm} (y - m) \quad \leftarrow \text{adjusted response}$$

$$= X\beta + \underbrace{\epsilon}_{\text{error term, need to find its variance}}$$

$$\text{Var}(\epsilon) = \left(\frac{d\eta}{dm} \right)^2 \underbrace{\text{Var}(y)}_{\propto \left(\frac{dm}{dn} \right)^2 v(n)}$$

(assuming $a(\phi)$ is a known constant)

$$W^{-1} = \left(\frac{d\eta}{dm} \right)^2 v(n) \quad a(\phi) \quad \left. \begin{array}{l} \text{needs to be} \\ \text{estimated} \end{array} \right\}$$

Now can estimate β :

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

Empirically fitting β_0 (logistic regression)

$$\log \frac{\pi_i}{1-\pi_i} = X\beta$$

$$(m_i, y_i) \rightarrow \hat{\pi}_i = \frac{y_i}{n_i}$$

$$\hat{y}_i = \log \frac{\hat{\pi}_i}{1-\hat{\pi}_i}$$



Fit a line and use its slope as β_0 then iterate.

Feb 16

(2) Fisher Scoring

log-likelihood for a single obs:

$$l = \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi)$$

$$n = \sum \beta_j x_j$$

$$\Rightarrow \frac{\partial l}{\partial \beta_j} = \frac{\partial l}{\partial \theta} \cdot \frac{d\theta}{d\mu} \cdot \frac{d\mu}{dn} \frac{\partial n}{\partial \beta_j}$$

$$= \frac{y - b'(\theta)}{a(\phi)} \cdot \frac{1}{V} \cdot \frac{\partial \mu}{\partial n} \cdot x_j$$

$$\mu = b'(\theta) \quad \frac{d\mu}{d\theta} = b''(\theta) = V(\mu)$$

$$W^{-1} = \left(\frac{d\mu}{d\theta} \right)^2 V(\mu) \Rightarrow \frac{d\mu}{dn} = W \frac{d\mu}{d\theta} V$$

$$\Rightarrow \frac{\partial l}{\partial \beta_j} = \frac{y - \mu}{a(\phi)} \frac{1}{V} W \frac{d\mu}{d\theta} V x_j$$

$$= \frac{W(y - \mu)}{a(\phi)} \frac{d\mu}{d\theta} x_j$$

If $a(\phi) = \phi$:

$$\frac{\partial l}{\partial \beta_j} = \frac{W(y - \mu)}{\phi} \frac{d\mu}{d\theta} x_j$$

If $a(\phi) = \phi/w_i$: y_1, \dots, y_n

$$\frac{\partial l}{\partial \beta_j} = \frac{w'(y - m)}{\phi} \frac{d n}{d m} x_j \quad \text{where } w_i \text{'s are absorbed into } w'$$

$$= \frac{\partial \sum l_i}{\partial \beta_j} = \sum_i \frac{\partial l_i}{\partial \beta_j}$$

$$= \sum_i \frac{w_i (y_{ij} - m_i) \frac{d n_i}{d m_i} x_{ij}}{\phi} = 0 \quad \left. \begin{array}{l} \text{solve to} \\ \text{get MLE} \end{array} \right\}$$

$$\frac{\partial^2 l}{\partial \beta_j \partial \beta_s} = w_i \frac{d n_i}{d m_i} x_{ij} \frac{\partial (y_{ij} - m_i)}{\partial \beta_s} + \frac{\partial w_i}{\partial \beta_s} \frac{d n_i}{d m_i} x_{ij} (y_{ij} - m_i)$$

$$\phi$$

taking expectation:

$$E \frac{\partial^2 l}{\partial \beta_j \partial \beta_s} = \frac{w_i \frac{d n_i}{d m_i} x_{ij} \frac{\partial (y_{ij} - m_i)}{\partial \beta_s} + 0}{\phi}$$

because $E(y_{ij} - m_i) = 0$

$$E \frac{\partial (y_{ij} - m_i)}{\partial \beta_s} = - \frac{\partial m_i}{\partial \beta_s}$$

$$= - \frac{w_i}{\phi} \frac{d n_i}{d m_i} x_{ij} \frac{d m_i}{d n_i} \frac{\partial n_i}{\partial \beta_s} = - \frac{w_i}{\phi} x_{ij} x_{is}$$

second deriv. in matrix form:

$$\frac{(X^T W X)}{\phi}, \quad W = \begin{bmatrix} w_1 & \dots & 0 \\ 0 & \ddots & \dots \\ 0 & \dots & w_n \end{bmatrix}_{n \times n} \quad X \in \mathbb{R}^{n \times p}$$

$$\Rightarrow \frac{(X^T W X)}{\phi} \in \mathbb{R}^{p \times p}$$

$$X^T W \left[(y - \mu) \frac{d\eta}{d\mu} \right] / \phi$$

Newton Raphson:

$$\beta_1 = \beta_0 + (X^T W X)^{-1} X^T W \left[(y - \mu) \frac{d\eta}{d\mu} \right]$$



phi's from these terms cancel out

$$= (X^T W X)^{-1} X^T W \left[X \beta_0 + (y - \mu) \frac{d\eta}{d\mu} \right]$$

Measuring the goodness of fit

1. $\ell(\hat{\mu}, \phi; y)$ - log-likelihood of the current model can be evaluated at the MLE. (maximized over β)

$\ell(y, \phi; \mu)$ - log likelihood of the fully saturated model w/ n parameters

- Derive canonical params from both models:

$$\hat{\theta} = \theta(\hat{m}) \text{ (current)}$$

$$\tilde{\theta} = \theta(y) \text{ (fully saturated)}$$

$$2 \left[l(\vec{y}, \phi; \vec{\tilde{\theta}}) - l(\vec{m}, \phi; \vec{\tilde{\theta}}) \right]$$

$$= \sum_i w_i \frac{\{ y_i (\tilde{\theta}_i - \hat{\theta}_i) - b(\tilde{\theta}_i) + b(\hat{\theta}_i) \}}{\phi}$$

$$= \frac{D(y, \hat{m})}{\phi} \quad \left. \right\} \text{scaled deviance}$$

$$\text{Deviance} \triangleq D(y, \hat{m})$$

$$\text{Normal } D = \sum_i (y_i - \hat{m}_i)^2 \text{ (exercise)}$$

$$\text{Poisson } D = \sum_i \left\{ y_i \log \frac{y_i}{\hat{m}_i} - (y_i - \hat{m}_i) \right\}$$

• Generalized Pearson χ^2 -statistic:

$$\chi^2 = \sum_i \frac{(y_i - \hat{m}_i)^2}{V(\hat{m}_i)}$$

• Residuals:

(1) Pearson's residual:

$$r_p = \frac{y - m}{\sqrt{V(m)}}$$

(2) Deviance residual:

$$D = \sum_i d_i$$

$$r_D = \underbrace{\text{sign}(y - m) \sqrt{d_i}}_{\text{deviance residual}}$$

\Rightarrow deviance is also a sum of squared residuals.

Log-linear model

- Poisson distribution:

$$P(Y=y) = \frac{\mu^y e^{-\mu}}{y!}$$

(a) cumulant generating function:

$$\begin{aligned} & \mu(e^t - 1) \\ &= \mu \left(t + \frac{t^2}{2!} + \frac{t^3}{3!} + \dots \right) \end{aligned}$$

All cumulants are equal to μ for Poisson

$$\frac{Y - \mu}{\sqrt{\mu}} \xrightarrow{\mu \uparrow} N(0, 1)$$

$$N: \quad 0 \quad 1 \quad 0 \quad 0$$

$$Y: \quad \mu \quad \mu \quad \mu \quad \mu$$

$$\frac{Y - \mu}{\sqrt{\mu}}: \quad 0 \quad 1 \quad O(\frac{\mu}{\sqrt{\mu}}) \quad O(\frac{\mu}{\mu^{4/2}}) \quad \text{remaining terms converge}$$

$$O(\mu^{-1/2}) \quad O(\mu^{-1}) \dots \text{to } 0$$

so $\frac{Y-\mu}{\sqrt{\mu}}$ cumulants converge to that of Standard Normal.

\sqrt{Y} is a variance-stabilizing transform for Poisson and is often used over \log .

$$g(Y) = g(\mu) + g'(\mu)(Y-\mu) + O(Y-\mu)$$

$$\text{Var}(g(Y)) = [g'(\mu)]^2 \text{Var } Y$$

$$= \left[\frac{1}{2} \mu^{-1/2} \right]^2 \mu$$

$$\text{Var}(\sqrt{Y}) \approx \frac{1}{4} \quad \left. \right\} \text{no longer depends on } \mu$$

Feb 21

log-linear model

assumption: $y_i \sim \text{Pois}(m_i)$

$$\log m_i = x_i^\top \beta \Rightarrow m_i = \exp(x_i^\top \beta)$$

$$\text{log-likelihood: } \Rightarrow \frac{\partial m_i}{\partial \beta} = e^{x_i^\top \beta} = m_i$$

$$l(\mu, y) = \sum_i (y_i \log m_i - m_i) + \text{const}$$

$$= \sum_i (y_i x_i^\top \beta - e^{x_i^\top \beta}) + \text{const}$$

$$\frac{\partial l}{\partial \beta_r} = \sum_i \frac{\partial l}{\partial m_i} \frac{\partial m_i}{\partial \beta_r} \frac{\partial \eta_i}{\partial \beta_r}$$

$$= \sum_i \left(\frac{y_i}{m_i} - 1 \right) m_i x_{ir}$$

$$= \sum_i (y_i - m_i) x_{ir}$$

$$\frac{\partial^2 l}{\partial \beta_r \partial \beta_s} = \sum_i -\frac{\partial m_i}{\partial \beta_s} x_{ir} = \sum_i -\frac{\partial m_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_r} x_{ir}$$
$$= -\sum_i m_i x_{is} x_{ir}$$

$$E\left[\frac{\partial^2 l}{\partial \beta_r \partial \beta_s}\right] = \sum_i m_i x_{is} x_{ir}$$

$$\Downarrow X^T W X, \quad W = \text{diag}(w_i)$$

$$\frac{\partial \ell}{\partial \beta} = X^T (y - w)$$

$$-E\left[\frac{\partial^2 \ell}{\partial \beta_r \partial \beta_s}\right] = X^T W X$$

Fisher scoring:

$$\beta_0 \quad \beta_1 = \beta_0 + (X^T W X)^{-1} X^T (y - w)$$

Deviance:

$$\begin{aligned} D(y, w) &= 2 \left[\ell(y, y) - \ell(\hat{w}, y) \right] \\ &= 2 \sum \left\{ y_i \log \frac{y_i}{\hat{w}_i} - (y_i - \hat{w}_i) \right\} \end{aligned}$$

Deviance residuals:

$$\text{sign}(y_i - \hat{w}_i) \sqrt{2 \{ y_i \log \frac{y_i}{\hat{w}_i} - (y_i - \hat{w}_i) \}}$$

Overdispersion:

$$\text{ex: } Y = \underbrace{Z_1 + Z_2 + \dots + Z_N}_{\text{indep. Poisson RVs}}$$

$Y \sim \text{Pois}$ if N fixed

If N is also an RV indep. of Z ,

Assume $EZ = \mu$ for all Z and $\text{var}(Z) = \sigma^2$

$$\begin{aligned} EY &= E[E(Y|N)] = E[NEZ] \\ &= ENEZ = \mu EN \end{aligned}$$

$$\begin{aligned} \text{var}Y &= E[\text{var}(Y|N)] + \text{var}[E(Y|N)] \\ &= E[N\text{var}Z] + \text{var}[NEZ] \\ &= (EN)\text{var}Z + \text{var}(N)(EZ)^2 \\ &= (EN)\mu + \text{var}(N)\mu^2 \end{aligned}$$

If $N \sim \text{Pois}$:

$$\text{var}Y = (EN)(\mu + \mu^2)$$

But $\text{var}Y > EY$ b/c N is random, additional layer of variability leads to overdispersion

ex: $Y_i | Z_i \sim \text{Pois}(Z_i)$, $Z_i \sim \text{Gamma}(\alpha, \beta)$

$$P(Y_i) = \int P(Y_i, Z_i) dZ_i$$

$P(Y_i)$ is not Poisson

Should be negative binomial. overdispersion.

$$\text{Var}(y) = \sigma^2 m, E(y) = m$$

$$\hat{\beta} =$$

$$D \sim \sigma^2 \chi^2$$

$\sim \sigma^2 m$

$$\frac{1}{n-p} \sum_i \frac{(y_i - m_i)^2}{m_i} = \hat{\sigma}^2 \quad \begin{matrix} \text{estimate of} \\ \text{dispersion param} \end{matrix}$$

log linear and multinomial response models

$$y_i \sim \text{Pois}(m_i), \quad \log(m_i) = \alpha + x_i^\top \beta = n_i$$

$$f(y_1, \dots, y_n) = f(y_+) f(y_1, \dots, y_n | y_+)$$

$$y_+ = \sum_i y_i \sim \text{Pois}(m_+), \quad m_+ = \sum_i m_i$$

$$f(y_+ | y_+) = \frac{f(y_+)}{f(y_+)} = \frac{\prod_{i=1}^n \frac{m_i^{y_i} e^{-m_i}}{y_i!}}{\frac{\sum_{i=1}^n m_i^{y_i} e^{-m_i}}{y_+!}}$$

$$= \frac{y_+!}{\prod_i y_i!} \prod_i \left(\frac{m_i}{m_+} \right)^{y_i} \frac{e^{-\sum m_i}}{e^{-m_+}}$$

$$\frac{m_i}{m_t} = \frac{e^{\alpha + x_i^\top \beta}}{\sum_i e^{\alpha + x_i^\top \beta}} = \frac{e^{x_i^\top \beta}}{\sum_i e^{x_i^\top \beta}}$$

$$\log f(y|\alpha, \beta) = \underbrace{\log(f(y_+))}_{\text{Poisson}} + \underbrace{\log(f(y|y_+))}_{\text{multinomial}}$$

Poisson

α, β

multinomial

β

plug in estimated

β from multinomial

first maximize this

estimate to get α estimate

y_+ is a sufficient statistic for α

$\pi_{11} \cdot \pi_{11}$	
π_{11}	π_{12}

T

π_{11}
 π_{12}
 assume T is fixed

Poisson way:

m_{11}	m_{12}
m_{21}	m_{22}

assume each cell's counts follows Poisson

ex: Independence

$$\text{multinomial: } \pi_{ij} = \pi_i + \pi_{+j}$$

$$\text{log linear: } m_{ij} = \exp\{\lambda + \lambda_i^x + \lambda_j^y\}$$

$$\log m_{ij} = \lambda + \lambda_i^x + \lambda_j^y \quad (*)$$

$$\pi_{ij} = \frac{m_{ij}}{M++} = \frac{e^{\lambda + \lambda_i^x + \lambda_j^y}}{\sum_i \sum_j e^{\lambda + \lambda_i^x + \lambda_j^y}}$$

$$\pi_{it} \leftarrow \frac{e^{\lambda_i^x} e^{\lambda_j^y}}{\sum_i e^{\lambda_i^x} \sum_j e^{\lambda_j^y}} \rightarrow \pi_{+j}$$

$\underbrace{\phantom{\sum_i e^{\lambda_i^x}}}_{\substack{\text{row factors} \\ \text{effects}}}$ $\underbrace{\phantom{\sum_j e^{\lambda_j^y}}}_{\substack{\text{column factors} \\ \text{effects}}}$

$$= \pi_i + \pi_{+j}$$

We usually constrain the params: $\sum_i \lambda_i^x = 1$

Deviance test tests for independence.

For an $I \times I$ table:

For identifiability, assume

$$\sum_i \lambda_i^x = 0 \quad \sum_j \lambda_j^y = 0$$

or

$$\lambda_I^x = 0 \quad \lambda_I^y = 0$$

params:

$$\lambda_1 + \lambda_{I-1}^* + \lambda_I^y = 2I - 1$$

I^2 data points

$$df: I^2 - (2I - 1) = (I - 1)^2$$

ex:

Residence in 1985

	NE	NW	S	W
NE	11607	100	306	124
residence in NW 1980	87	13677	515	302
S	172	255	17819	270
W	63	176	286	10192

under independence:

$$\log \mu_{ij} = \lambda + \lambda_i^X + \lambda_j^Y$$

to test for independence:

$$\log \mu_{ij} = \lambda + \lambda_i^X + \lambda_j^Y + \delta_i I(i=j)$$



Quasi-Independence model

$$I \subset QI$$

 independence model

Feb 23

$$y_i \sim \text{Pois}_+(\mu_i)$$

$$\log \mu_i = x_i^T \beta \Rightarrow \mu_i = e^{\alpha + x_i^T \beta}, \mu_+ = \sum e^{\alpha + x_i^T \beta}$$

$$P(y_1, \dots, y_n) = \prod_{i=1}^n \frac{\mu_i^{y_i} e^{-\mu_i}}{y_i!}$$

$$\begin{aligned} \mu_+ &= \sum_i \mu_i \\ &= \frac{y_+ (\mu_+)^{y_+}}{y_+ (\mu_+)^{y_+}} \prod_{i=1}^n \frac{\mu_i^{y_i} e^{-\mu_i}}{y_i!} \\ &= \underbrace{\frac{(\mu_+)^{y_+} e^{-\mu_+}}{y_+!}}_{\text{Poisson}} \underbrace{\prod_i \left(\frac{\mu_i}{\mu_+} \right)^{y_i}}_{\text{multinomial}} \end{aligned}$$

$$\ell(\alpha, \beta, y) = y_+ \log \mu_+ - \mu_+ + \ell(\beta | y)$$

$$\frac{\mu_i}{\mu_+} = \frac{e^{x_i^T \beta}}{\sum e^{x_i^T \beta}}$$

reparameterize the log-likelihood

$$\ell(\mu_+, \beta; y)$$

$$\hat{\mu}_+ = y_+$$

Independence model:

$$\log \pi_{ij} = \lambda + \lambda_i^x + \lambda_j^y$$

Quasi Independence model:

$$\log \pi_{ij} = \lambda + \lambda_i^x + \lambda_j^y + \delta_i I(i=j)$$

Symmetry

$$\pi_{ij} = \pi_{ji} \quad \forall i, j$$

$$\log \pi_{ij} = \lambda + \lambda_i^x + \lambda_j^y + \lambda_{ij}$$

$$\text{where } \lambda_{ij} = \lambda_{ji}$$

$$\log \pi_{ji} = \lambda + \lambda_j^y + \lambda_i^x + \lambda_{ji}$$

Don't have x, y , rows and columns share the same type of effects

Symmetry \Rightarrow marginal homogeneity (MH)

$$\sum \text{marginal row sums} = \sum \text{marginal col sums}$$

$$\pi_{rj} = \sum_i \pi_{ij} = \sum_i \pi_{ji} = \pi_{rj} +$$

MH $\not\Rightarrow$ symmetry

$2 \times 2: MH \Rightarrow S$

π_{11}	π_{12}
π_{21}	π_{22}

$$\begin{aligned}\pi_{11} + \pi_{12} &= \pi_{11} + \pi_{21} \\ \Rightarrow \pi_{12} &= \pi_{21}\end{aligned}$$

$3 \times 3:$

x	y	z
z	x	y
y	z	x

MH but not S

Free params:

$$\log \mu_{ij} = \lambda + \underbrace{\lambda_i}_{\frac{I-1}{2}} + \underbrace{\lambda_j}_{\frac{I(I-1)}{2}} + \lambda_{ij} \quad \left. \right)$$

Only need to consider
off diagonals in
one half b/c of
symmetry

$$df = I - (I-1) - \frac{I(I-1)}{2}$$

$$= \frac{I(I-1)}{2}$$

Quasi-symmetry: $S \subset QS$

$$\log M_{ij} = \lambda + \lambda_i^x + \lambda_j^y + \frac{\lambda_{ij}}{z}$$

models overall mean where $\lambda_{ij} = \lambda_{ji}$

params: $1 + I - 1 + I - 1 + \frac{I(I-1)}{2}$

$$\begin{aligned} df &= I^2 - (\# \text{params}) \\ &= \frac{(I-1)(I-2)}{2} \end{aligned}$$

$QS + MH \Rightarrow S :$

consider multinomial probabilities implied by QS .

$$\log M_{ij} = \lambda + \lambda_i^x + \lambda_j^y + \lambda_{ij}$$

$$= \log M_{++} + \log \frac{M_{ij}}{M_{++}}$$

$$\approx \log M_{++} + \log \pi_{ij}$$

$$\pi_{ij} \propto e^{\lambda_i^x} e^{\lambda_j^y} e^{\lambda_{ij}}$$

$$\pi_{ab} \propto \alpha_a \beta_b \gamma_{ab}$$

| | interaction
 ath row bth column effects

$$\begin{aligned}
 &= \frac{\alpha_a}{\beta_a} \beta_a \beta_b \gamma_{ab} \\
 &= l_a \delta_{ab} \quad \leftarrow \delta_{ab} = \delta_{ba} \text{ (symmetric quantity)}
 \end{aligned}$$

$$\pi_{ab} = l_a \delta_{ab}$$

$$\pi_{ab} \stackrel{?}{=} \pi_{ba} \quad \text{we know: } \pi_{ji} = \pi_{ij} \quad (\text{MH})$$

$$\pi_{ji} = \sum_i \pi_{ji} \quad \pi_{ij} = \sum_i \pi_{ij}$$

$$= \sum_i l_j \delta_{ji} \quad = l_i \sum_i \delta_{ij}$$

$$= l_j \sum_i \delta_{ji} \quad = \sum_i l_i \delta_{ij}$$

$$= l_j \sum_i \delta_{ij}$$

$$l_j = \sum_i l_i \frac{d_{ij}}{\sum_i d_{ij}}$$

$$= \sum_i l_i w_i$$

$$w_i \triangleq \frac{d_{ij}}{\sum_i d_{ij}}$$

$l_1 \leq l_2 \leq \dots \leq l_I$ and $l_1 < l_I$ if they are not all the same.

$$l_1 = \sum_i l_i w_i \geq l_1 \sum_i^{I-1} w_i + l_I w_I$$

$$\geq l_1 \sum_{i=1}^{I-1} w_i \\ = l_1$$

which is a contradiction. So

$$l_1 = l_2 = \dots = l_I$$

$$\Rightarrow \pi_{ab} = l_a f_{ab} = l_b f_{ba} = \pi_{ba}$$

Quasi-likelihood

$$E(y) \sim X\beta$$

$$\text{var}(y) \sim V(\mu)$$

(1) Suppose response variable y has

mean: μ

variance: $\sigma^2 V(\mu)$

$$U = \frac{y - \mu^{(\beta)}}{\sigma^2 V(\mu)} = 0$$

behaves like devia.
of log-likelihood.

Solve for estimate of β .

$$E_y l'(\theta; y) = 0$$

$$-E l''(\theta; y) = \text{var}(l'(\theta; y))$$

$$= E[l'(\theta; y)]^2$$

$$\cong \frac{1}{\sigma^2 V(\mu)}$$

) asymptotically

$$E(U) = 0 \quad (\text{only } y \text{ is random})$$

$$\text{Var}(U) = \frac{\sigma^2 V(\mu)}{(\sigma^2 V(\mu))^2} = \frac{1}{\sigma^2 V(\mu)}$$

$$-E \frac{\partial U}{\partial \mu} = E \left[\frac{-1}{\sigma^2 V(\mu)} + (y-\mu) \cancel{\left(\quad \right)} \right]$$

$$= \frac{1}{\sigma^2 V(\mu)}$$

\uparrow
 $\text{var}(y)^{-1}$

So U behaves like deriv. of log-likelihood.

March 2

Quasi-Likelihood

$\mu, V(\mu)$

$$\frac{y - \mu}{\sigma^2 V(\mu)}$$

Ex 2:

$$V(m) = m, \sigma^2 = 1$$

$$\begin{aligned} Q(m, y) &= \int_y^m \frac{y-t}{t} dt \\ &= (y \log t - t) \Big|_y^m \\ &= (y \log m - m) - (y \log y - y) \end{aligned}$$

Maximising the Quasi-likelihood:

• Differentiate $Q(m, y)$ wrt β

$$\begin{aligned} U(\beta)_{px1} &= \frac{\partial m^T}{\partial \beta} \frac{V^{-1}(y-m)}{\sigma^2} \\ &= D_{pxn}^T \frac{V^{-1}(y-m)}{\sigma^2} \stackrel{\text{set}}{=} 0 \end{aligned}$$

$$\text{where } D = (D_{ir})_{n \times p} \triangleq \left(\frac{\partial m_i}{\partial \beta_r} \right)_{n \times p}$$

$$\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \sim \begin{bmatrix} m_1 \\ \vdots \\ m_n \end{bmatrix} \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

$$\begin{bmatrix} V(m_1) & & 0 \\ & \ddots & \\ 0 & & V(m_n) \end{bmatrix}_{n \times n} \quad \begin{bmatrix} y_1 - m_1 \\ \vdots \\ y_n - m_n \end{bmatrix}$$

$$\begin{aligned}
 & -E \frac{\partial u(\beta)}{\partial \beta} \\
 &= -E \left[-\frac{D^T V^{-1} D}{\sigma^2} + \frac{\partial D^T V^{-1}}{\partial \beta} \frac{(y - m)}{\sigma^2} \right] \\
 &= \frac{D^T V^{-1} D}{\sigma^2} \stackrel{\Delta}{=} i_{\beta} \quad \text{znd deriv of Quasi likelihood}
 \end{aligned}$$

0 in expectation

Newton-Raphson:

✓ dispersion param
 σ^2 gets cancelled

$$\begin{aligned}
 \beta_1 &= \beta_0 + \left(\frac{D_0^T V_0^{-1} D_0}{\sigma^2} \right)^{-1} \frac{D_0^T V_0^{-1} (y - m_0)}{\sigma^2} \\
 &= \beta_0 + (D_0^T V_0^{-1} D_0)^{-1} D_0^T V_0^{-1} (y - m_0)
 \end{aligned}$$

$$\text{var}(\hat{\beta}) \simeq \sigma^2 (D^T V^{-1} D)^{-1}$$

$$\hat{\sigma}^2 = \frac{1}{n-p} \sum_{i=1}^n \frac{(y_i - m_i)^2}{V(m_i)} = \frac{\chi^2}{n-p}$$

Optimal estimating functions

Quasi-score $U(\beta) = D^T V^{-1} (y - m) / \sigma^2$ is a special case of estimating function.

(1) Estimating function $g(y; \theta)$: a function of data y and parameter θ with zero mean for all parameter values $E_y g(y; \theta) = 0 + \theta$

$$H_{p \times n}^T (y - m)_{n \times 1}, \quad H(\beta)_{n \times p}$$

$\underbrace{\quad}_{p \times 1}$. Get p estimating functions, set to 0, get p estimating eqns.

$U(\beta) = D^T V^{-1} (y - m) / \sigma^2$ is the optimal combination of the elementary est function within the class of linear estimating function.

$$(y - m) \rightarrow h(y; \beta) \triangleq H^T (y - m) = 0 \Rightarrow \hat{\beta} \Rightarrow \hat{\beta}$$

Asymptotically, all linear functions of $\hat{\beta}$ have variances at least as big as the variance of $\hat{\beta}$.

$$\text{i.e. } \text{var}(\tilde{\beta}) \geq \text{var}(\hat{\beta}).$$

$$h(y; \hat{\beta}) = h(y; \beta) + h'(y; \beta)(\hat{\beta} - \beta) + o(\hat{\beta} - \beta)$$

↑
can replace w/ its expectation

$$\text{then } \hat{\beta} - \beta \simeq -[E h'(y; \beta)]^{-1} h(y; \beta)$$

$$\begin{aligned} &= \left[-E(-H^T D + \frac{\partial H}{\partial \beta}(y - \mu)) \right]^{-1} h(y; \beta) \\ &= (H^T D)^{-1} h(y; \beta) \\ &= (H^T D)^{-1} H^T (y - \mu) \end{aligned}$$

$$\begin{aligned} \text{var}(\hat{\beta}) &\simeq \text{var}(H^T D)^{-1} H^T (y - \mu) \\ &= (H^T D)^{-1} H^T V H (D^T H)^{-1} \sigma^2 \end{aligned}$$

$$\text{var}(\hat{\beta}) = \sigma^2 (D^T V^{-1} D)^{-1}$$

$$\text{var}(\hat{\beta}) - \text{var}(\hat{\beta}) \geq 0$$

non-negative definite matrix

Need to show $(H^T D)^{-1} H^T V H (D^T H)^{-1} (D^T V^{-1} D)^{-1} \geq 0$

$$\Leftrightarrow b/c A - B \geq 0 \Leftrightarrow B^{-1} - A^{-1} \geq 0$$

$$D^T V^{-1} D - (D^T H) (H^T V H)^{-1} H^T D \geq 0$$

$$\Leftrightarrow$$

$$D^T [V^{-1} - H (H^T V H)^{-1} H^T] D \geq 0$$

$$\Leftrightarrow$$

$$V^{-1} - H (H^T V H)^{-1} H^T \geq 0$$

$$V^{1/2} [V^{-1} - H (H^T V H)^{-1} H^T] V^{1/2} \geq 0$$

$$\Leftrightarrow$$

$$I - V^{1/2} H (H^T V H)^{-1} H^T V^{1/2} \geq 0$$

$$\tilde{X} \triangleq V^{1/2} H$$

$$I - \tilde{X} (\tilde{X}^T \tilde{X})^{-1} \tilde{X} \geq 0$$

March 7

Optimal estimating eqns

$$\begin{bmatrix} D' V^{-1} (Y - \mu) \\ \begin{bmatrix} 0 & \dots & 0 \end{bmatrix}^{-1} \begin{matrix} y_1 - m_1 \\ y_2 - m_2 \\ \vdots \\ y_n - m_n \end{matrix} \end{bmatrix} = E\{g_i(Y, \theta) | A_i\} = 0$$

Time ordered sequence:

$$Y_1, Y_2, \dots, Y_n$$

$$Y_{(t)} = (Y_1, \dots, Y_t)$$

$$E\{g_t(Y_t, \theta) | Y_{t-1}\} = 0$$

$$V \triangleq \text{diag}\{\text{Var}(g_i | A_i)\}$$

$$\text{Dir} \triangleq -E\left\{\frac{\partial g_i(Y, \theta)}{\partial \theta} \mid A_i\right\}$$

$$U(\theta, y) = D' V^{-1} g$$

Ex: Y 's are generated by an autoregressive process.

$$Y_t = \theta Y_{t-1} + \varepsilon_t$$

$$Y_0 = \varepsilon_0$$

$$\text{where } \varepsilon_t \stackrel{iid}{\sim} N(0, 1)$$

Want to estimate θ :

Need a function that is 0 in conditional expectation.

$$\mathbb{E} \{ y_t - \theta y_{t-1} \mid Y_{t-1} \} = 0$$

$$g_t = y_t - \theta y_{t-1}$$

$$g^* = \frac{y_t}{\theta} - y_{t-1} \quad (\text{if } \theta \neq 0)$$

$$g^* = Bg, \quad U(\theta, y) = D' V^{-1} g$$

$$\text{Var}(g^*) = B \text{Var}(g) B' = V^*$$

$$= BV B'$$

$$\frac{\partial g^*}{\partial \theta} = B \frac{\partial g}{\partial \theta} = BD = D^*$$

$$D^* V^{*-1} g^*$$

$$= (BD)' (BV B')^{-1} Bg$$

$$= D' B' (B')^{-1} V^{-1} B^{-1} Bg$$

$$Dir = -E\left(\frac{\partial g_i}{\partial \theta_i} \mid A_i\right)$$

$$= D' V^{-1} g$$

$$D = \begin{bmatrix} y_0 \\ y_1 \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} y_0 & y_1 & \dots \end{bmatrix} \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}^{-1} \begin{bmatrix} y_1 - \theta y_0 \\ y_2 - \theta y_1 \\ \vdots \\ y_n - \theta y_{n-1} \end{bmatrix} = \sum_t y_{t-1} g_t$$

$$= \sum_t y_{t-1} (y_t - \theta y_{t-1}) = 0$$

$$\hat{\theta} = \frac{\sum_t y_{t-1} y_t}{\sum_t y_{t-1}^2}$$

$$\begin{aligned} g^*: V(\theta, y) &= \sum_t E\left(\frac{y_t}{\theta^2} | y_{t-1}\right) \theta^2 \left(\frac{y_t}{\theta} - y_{t-1}\right) \\ &= \sum_t \frac{\theta y_{t-1}}{\theta^2} \theta^2 \left(\frac{y_t}{\theta} - y_{t-1}\right) \\ &= \sum_t y_{t-1} (y_t - \theta y_{t-1}) \end{aligned}$$

MLE:

$$\prod_t \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(y_t - \theta y_{t-1})^2}{2}\right\}$$

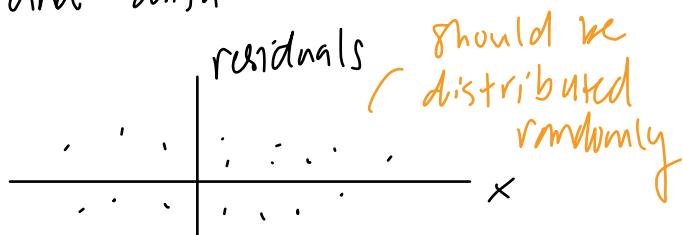
$$l = \sum_t \frac{-(y_t - \theta y_{t-1})^2}{2} + c$$

$$\frac{\partial l}{\partial \theta} = \sum_t \frac{2(y_t - \theta y_{t-1}) y_{t-1}}{2} = 0$$

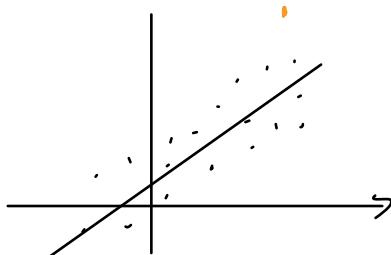
\Rightarrow optimal estimating eqns and MLE return the same estimate.

Model Diagnosis

- 1) Check systematic discrepancies between model and data



- 2) Check isolated discrepancies



{ Formal: tests
Informal: plots

asymptotically equivalent

(1) Drop-in-deviance
(likelihood ratio) test

M_0 : p params

M_1 : p+k params

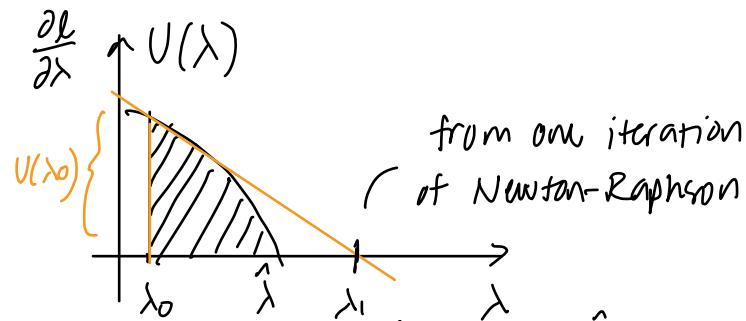
$$D = 2(\log L_1 - \log L_0)$$

$$= 2(l_1 - l_0)$$

(2) Score test

$$S(\lambda_0) = U(\lambda_0)^T_i U^{-1}(\lambda_0) U(\lambda_0)$$

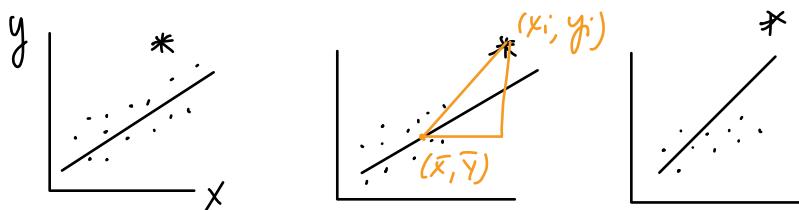
Inv. of Fisher
eval. at MLE



$$l(\hat{\lambda}) - l(\lambda_0) = \int_{\lambda_0}^{\hat{\lambda}} \frac{d l}{d \lambda} d \lambda = \int_{\lambda_0}^{\hat{\lambda}} U(\lambda) d \lambda$$

$$\frac{1}{2} S(\lambda_0) = \frac{1}{2} U(\lambda_0) (\lambda_1 - \lambda_0)$$

Isolated discrepancies:



$$\begin{aligned}
 y - \bar{y} &= \gamma \frac{s_y}{s_x} (x - \bar{x}) \\
 &= \frac{s_{xy}}{s_x^2} (x - \bar{x}) \\
 &= \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} (x_i - \bar{x})
 \end{aligned}$$

$$\beta = \frac{\sum_i (x_i - \bar{x})^2 \left(\frac{y_i - \bar{y}}{x_i - \bar{x}} \right)}{\sum_i (x_i - \bar{x})^2}$$

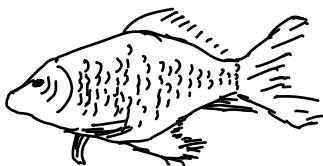
$$= \sum_i w_i \left(\frac{y_i - \bar{y}}{x_i - \bar{x}} \right)$$

linear combination of
slope contributed by each datapoint

if x is very different from \bar{x} , it will have a bigger weight.



cabbage



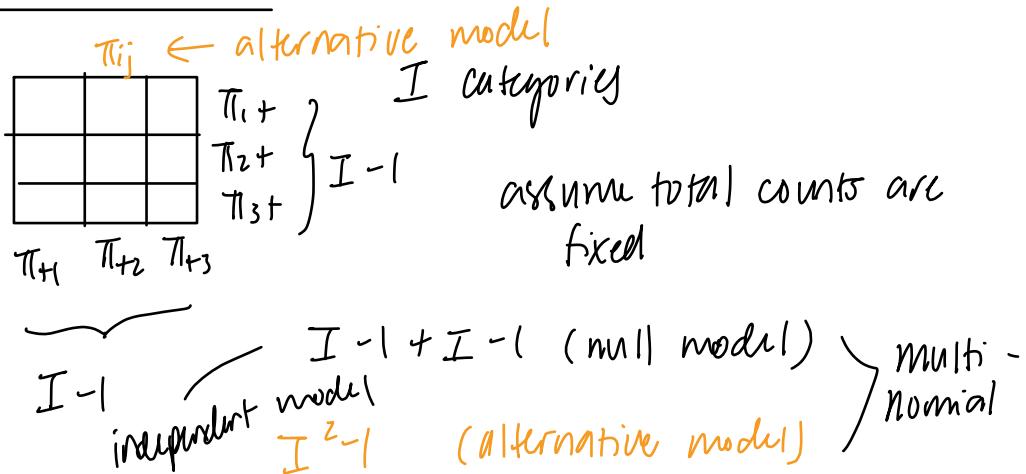
$$\hat{y} = X\hat{\beta}$$

$$H = X(X'X)^{-1}X'$$

$$\hat{y}' = Hy$$

$$H = \begin{bmatrix} \ddots & & \\ & h_{ii} & \\ \ddots & & \ddots \end{bmatrix}$$

March 7 04



difference:

$$I^2 - I - 2(I - 1)$$

$$= I^2 - 2I + 1$$

$$= (I - 1)^2$$

χ^2 will have this many d.f.

alternative model:

T

i, j

pick π_{11} to be baseline

introduce an indicator / indicators

$$\begin{bmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{II} \end{bmatrix} = \begin{bmatrix} 1 & \alpha_{11} & \alpha_{12} & \dots & \alpha_{I-1} \\ 1 & \alpha_{21} & \dots & & \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_{11} \\ \beta_{12} \\ \vdots \\ \beta_{I-1} \end{bmatrix} \quad I^2 - 1$$

exercisit: converting models to matrix form

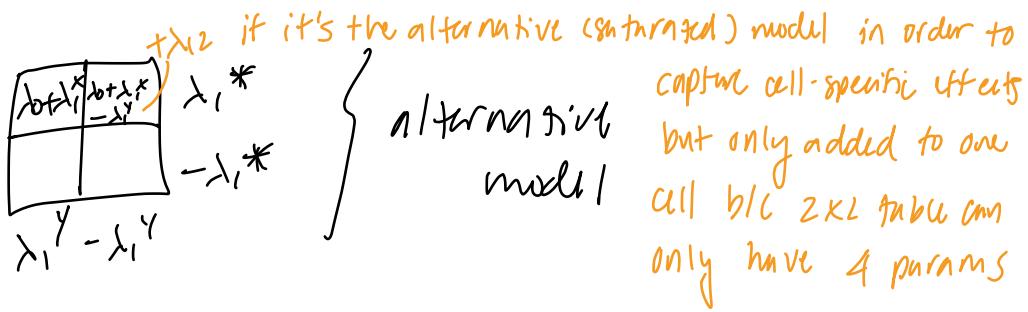
Poisson log-linear is the same:

$$\log m_{ij} = \lambda_0 + \lambda_i^x + \lambda_j^y \quad | \text{ in row } 2 \text{ or not} \quad | \text{ indicator of obs.} \quad | \text{ column indicators}$$

$$\begin{bmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{II} \end{bmatrix} \rightarrow \begin{bmatrix} \log m_{11} \\ \log m_{12} \\ \vdots \\ \log m_{II} \end{bmatrix} = \begin{bmatrix} 1 & r_{11} & r_{12} & c_{11} & c_{12} \\ \vdots & r_{21} & r_{22} & \vdots & \vdots \\ 1 & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \lambda_0 \\ \lambda_1^x \\ \lambda_2^x \\ \lambda_I^x \\ \lambda_1^y \\ \vdots \\ \lambda_I^y \end{bmatrix}$$

$2I - 1$ params

↪ same as independent multinomial model.



$$z_{ij} \sim N(\mu, \sigma^2)$$

$$z_{ij} = y_i + \varepsilon_{ij} \sim N(0, \tau^2)$$

$$\text{cov}(z_{ij}, z_{i'j'})$$

$$= \text{cov}(y_i + \varepsilon_{ij}, y_{i'} + \varepsilon_{i'j'}) \quad (i \neq i')$$

$$= 0$$

$$\text{cov}(z_{ij}, z_{in}) \quad j \neq n$$

$$= \text{cov}(y_i + \varepsilon_{ij}, y_i + \varepsilon_{in})$$

$$= \underbrace{\text{cov}(y_i, y_i)}_{\text{marginal covariance}} = \text{var}(y_i) = \sigma^2$$

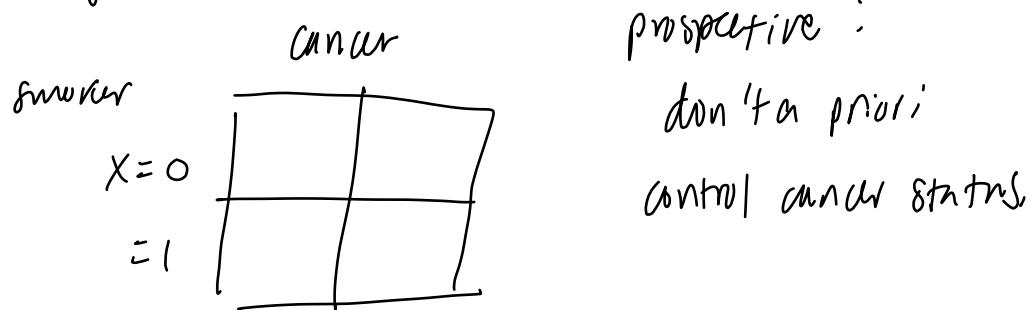
↓
non-zero
correlation

$$z_{ij} \mid y_i \stackrel{iid}{\sim} N(y_i, \tau^2), y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

$$\text{cov}(z_{ij}, z_{ik} | y_i) \neq \text{cov}(z_{ij}, z_{ik})$$

0

$$\text{logit } P(y=1|X) = \beta_0 + \beta_1 X$$



Conditional on being a smoker, what's
the prob. of getting cancer? \rightarrow same
as population prob.

March 9

$$r = y - \hat{y} = y - Hy \\ = (I - H)y$$

$$\text{var } r = (I - H)\sigma^2 I (I - H)$$

$$= (I - H)\sigma^2 I$$

$$\text{var} \begin{bmatrix} y_1 - \hat{y}_1 \\ \vdots \\ y_n - \hat{y}_n \end{bmatrix} \rightarrow \underbrace{\begin{bmatrix} & & \\ & 0 & \\ & & \end{bmatrix}}_{I - H} (I - H) \sigma^2$$

Raw materials for model checking

(1) Linear models

\hat{y} : fitted values

$r = y - \hat{y}$: residuals (raw)

s^2 : residual variance

$H = X(X'X)^{-1}X'$: hat matrix

$h_i \triangleq h_{ii}$: leverage $= x_i'(X'X)^{-1}x_i$
 $[\dots] \quad [:]$

other types of residuals:

standardized:

$$\frac{y_i - \hat{m}_i}{\sqrt{1-h_i}}$$

standardized:
standardized:

$$r_i' = \frac{y_i - \hat{m}_i}{s \sqrt{1-h_i}}$$

fitted w/ problematic
point

deleted:

$$r_i^* = \frac{y_i - \hat{m}_{(i)}}{s_{(i)} \sqrt{1+h_{(i)}}}$$

ith data point
removed

$$\hat{m}_{(i)} = \hat{x}_i' \hat{\beta}(i)$$

↑
obtained w/o
ith data point

$$h_{(i)} = \hat{x}_i' (\hat{X}_{(i)}' \hat{X}_{(i)})^{-1} \hat{x}_i$$

y_i and $\hat{m}_{(i)}$ are independent, but y_i and \hat{m}_i are not.
that's why variance in deleted residual is additive.

$$r_i^* = \frac{y_i - \hat{m}_{(i)}}{s_{(i)} \sqrt{1+h_{(i)}}} = \frac{y_i - \hat{m}_i}{s_{(i)} \sqrt{1-h_i}}$$

Pf: / wlog, consider obs. 1

$$\text{var}(y_1 - \hat{m}_{(1)}) = \text{var}(y_1 - \hat{x}_1' (\hat{X}_{(1)}' \hat{X}_{(1)})^{-1} \hat{X}_{(1)}' y_{(1)})$$

$$\text{var}(y_1 - \hat{m}_{(1)})$$

$$\text{var}(Ay) = A \text{var}(y) A'$$

$$= \text{var} \left(\begin{bmatrix} 1 & X_{(1)}' (X_{(1)}' X_{(1)})^{-1} X_{(1)}' \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \right)$$

$$= \left[1 + X_{(1)}' (X_{(1)}' X_{(1)})^{-1} X_{(1)}' \right] \sigma^2 I \begin{bmatrix} 1 \\ -X_{(1)} (X_{(1)}' X_{(1)})^{-1} X_{(1)}' \end{bmatrix}$$

$$= \left[1 + X_{(1)}' (X_{(1)}' X_{(1)})^{-1} X_{(1)}' X_{(1)} (X_{(1)}' X_{(1)})^{-1} X_{(1)}' \right] \sigma^2$$

$$= \left[1 + h_{(1)} \right] \sigma^2$$

$$= \left[1 + h_{(1)} \right] \sigma^2$$

$$\text{so } r_i^* = \frac{y_i - \hat{m}_{(i)}}{s(i) \sqrt{1+h_{(i)}}}$$

Next, show $h_{(i)} = \frac{h_i}{1-h_i}$, wlog, consider h_i :

$$X = \begin{pmatrix} X_1' \\ \vdots \\ X_n' \end{pmatrix} = \begin{pmatrix} X_1' \\ X_{(1)}' \end{pmatrix}$$

$$H = X(X'X)^{-1}X' = \begin{pmatrix} X_1' \\ X_{(1)}' \end{pmatrix} \left[(X_1 X_{(1)}) \begin{pmatrix} X_1' \\ X_{(1)}' \end{pmatrix} \right]^{-1} \begin{pmatrix} X_1' \\ X_{(1)}' \end{pmatrix}$$

$$= \begin{pmatrix} X_1' \\ X_{(1)}' \end{pmatrix} \left[X_1 X_1' + X_{(1)}' X_{(1)} \right]^{-1} (X_1 X_{(1)})'$$

$$(A - BD^{-1}C)^{-1} = A^{-1} + A^{-1}B(D - CA^{-1}B)^{-1}CA^{-1}$$

$$A: X_{(1)}' X_{(1)}, \quad B: X_1, \quad C: X_1', \quad D = I$$

$$\left[X_1 X_1' + X_{(1)}' X_{(1)} \right]^{-1}$$

$$= [X_{(1)}' X_{(1)}]^{-1} + [X_{(1)}' X_{(1)}]^{-1} X_1 (-1 - X_1' (X_{(1)}' X_{(1)})^{-1} X_1) X_1' (X_{(1)}' X_{(1)})^{-1}$$

$$= [X_{(1)}' X_{(1)}]^{-1} [X_{(1)}' X_{(1)}]^{-1} X_1 (1 + h_{(1)})^{-1} X_1' (X_{(1)}' X_{(1)})^{-1}$$

$$h_1 = X_1' [X_1 X_1' + X_{(1)}' X_{(1)}]^{-1} X_1$$

$$= X_1' [X_{(1)}' X_{(1)}]^{-1} X_1 - \frac{X_1' [X_{(1)}' X_{(1)}]^{-1} X_1 \cdot X_1' [X_{(1)}' X_{(1)}]^{-1} X_1}{1 + h_{(1)}}$$

$$h_1 = h_{(1)} - \frac{h_{(1)}^2}{1 + h_{(1)}} = \frac{h_{(1)}}{1 + h_{(1)}}$$

$$h_{(1)} = \frac{h_1}{1 - h_1}$$

$$y_1 - \hat{m}_1 = y_1 - x_1' (X'X)^{-1} X'y$$

$$= y_1 - x_1' \left[(x_1 x_{(1)'}') \begin{pmatrix} x_1 \\ x_{(1)'} \end{pmatrix} \right]^{-1} (x_1 x_{(1)'}') \begin{pmatrix} y_1 \\ y_{(1)} \end{pmatrix}$$

$$= (1-h_1)(y_1 - \hat{m}_{(1)})$$

$$\Rightarrow h_{(1)} = \frac{h_1}{1-h_1}$$

Show these steps
at home.

$$\frac{y_1 - \hat{m}_1}{\sqrt{1-h_1}} = \frac{y_1 - \hat{m}_{(1)}}{\sqrt{1+h_{(1)}}}$$

$$r_i^* = r_i' \frac{s_i}{s_{(i)}}$$

GLMs:

$$y \rightarrow z \text{ (adjusted response)}$$

$$\hat{m} \rightarrow \hat{\eta}$$

$$S^2 \rightarrow \hat{\phi}$$

$$H \rightarrow W'^2 X (X' W X) X' W'^2$$

$$W = V^{-1} \left(\frac{\partial \hat{m}}{\partial \hat{\eta}} \right)^2$$

$$r_i' = \frac{y_i - \hat{m}_i}{\sqrt{\hat{\phi} V(\hat{m}_i)(1-h_i)}}$$

$$g(y) \stackrel{\sim}{=} g(\mu) + g'(\mu)(y - \mu)$$

$$\tilde{z} = X\beta + \frac{\partial n}{\partial \mu}(y - \mu)$$

$$= X\beta + \varepsilon \rightarrow \text{var}(\varepsilon) = \left(\frac{\partial n}{\partial \mu} \right)^2 \phi V(\mu)$$

$$W^{-1} = \phi V(\mu) \left(\frac{\partial n}{\partial \mu} \right)^2$$

$$W^{1/2} \tilde{z} = W^{1/2} X\beta + W^{1/2} \varepsilon$$

$$\hat{z} = \tilde{X}\hat{\beta} + \tilde{\varepsilon} \quad \leftarrow \text{reduced to linear regression}$$

so

$$\underbrace{\tilde{X}(\tilde{X}'\tilde{X})^{-1}\tilde{X}'}_{\text{hat matrix}} = W^{1/2} X(X'W X)^{-1} X' W^{1/2}$$

hat matrix

$$\text{var}(\hat{\varepsilon}) = W^{1/2} \phi V W^{1/2} I$$

$$\frac{\hat{z}_i - \tilde{x}_i'\hat{\beta}}{\sqrt{1-h_i}\phi} = \phi I$$

$$\hat{\varepsilon} = \frac{W^{1/2} (z - X\beta)}{\sqrt{(1-h_i)\phi}}$$

$$= \frac{W_i^{1/2} \left(\frac{\partial m_i}{\partial h_i} \right) (y_i - \hat{m}_i)}{\sqrt{(1-h_i) \phi}}$$

$$W = \left(\frac{\partial m}{\partial h_i} \right)^2 V^{-1}$$

$$W^{1/2} = \left(\frac{\partial m}{\partial h_i} \right) V^{-1/2}$$

$$= \frac{\left(\frac{\partial m_i}{\partial h_i} \right) \left(\frac{\partial m_i}{\partial h_i} \right) V_i^{-1/2} (y_i - \hat{m}_i)}{\sqrt{(1-h_i) \phi}}$$

$$= \frac{(y_i - \hat{m}_i)}{\sqrt{\phi V(h_i)(1-h_i)}}$$

$$\begin{aligned} \sum_i h_i &= p = \text{tr}(H) \\ &= \text{tr}(X(X'X)^{-1}X') \\ &= \text{tr}\left((X'X)^{-1}X'X\right) \\ &= \text{tr}(I_{p \times p}) = p \end{aligned}$$

so each datapoint should contribute p/n

if $h_i \geq \frac{2p}{n}$, it's high leverage

in linear models, large h_i = covariate is different. But not necessarily in GLM.
 Because \hat{X} contains a weight. A very different datapoint can have a small weight and thus exert less influence.

Cook's distance:

Linear:

$$\hat{\beta}_{(i)} - \hat{\beta}$$

$$D_i = \frac{(\hat{\beta}_{(i)} - \hat{\beta})' (X'X) (\hat{\beta}_{(i)} - \hat{\beta})}{p s^2}$$

$$= \frac{r_i^2}{p} \cdot \frac{h_i}{1-h_i}$$

monotone function of leverage

\Rightarrow large leverage \rightarrow large D_i

GLM:

$$D_i = \frac{(\hat{\beta}_{(i)} - \hat{\beta})(X'W X)(\hat{\beta}_{(i)} - \hat{\beta})}{P \phi}$$

\uparrow
params

add weight matrix from
linear approx.

March 14

Matched aux-control design

Y : outcome { 1 disease
0 control }

X: primary risk factor: smoking status

secondary risk factors: alcohol

2: potential confounders : age , education levels, etc.

sample the data so that they match on potential confounders.

$i = 1 \quad x \quad x \quad x^{m_i}$ N cases, for each case all but
 \vdots
 N m_i controls matched w/
 potential confounders.

Fit a logistic regression model:

$$\text{logit } \pi = \beta_0 + x\beta_x + z\beta_z$$

Assuming linear dependence of logit on confounders (not guaranteed to be true in the data).

how to get rid of this assumption?

logit $\pi_i = \alpha_i + \tilde{x}_i^T \beta$ for the i^{th} matching
each matching group has its own intercept \tilde{x}_i^T primary risk factor

- b/c all datapoints in the group have the same covariates values.
- No longer making assumptions about confounders.
- Issue: Many parameters to fit.
 - N cases, N_K datapoints (K -constant)
 - $O(N)$ params
 - ⇒ Not enough datapoints to estimate in a consistent way
 - H^{-1} will be expensive to compute in Newton-Raphson
 - can't directly use this model

conditional likelihood approach

- Used to handle above issues.

$$P(Y_{i0} = 1, Y_{i1} = 0, Y_{i2} = 0, \dots, Y_{im_i} = 0) \sum_j Y_{ij} = 1$$

(show this doesn't depend on x_i)

- only one obs. is 1, rest are controls for i th matching group.

$$\left\{ \begin{array}{ll} Y_{i0} = 1 & \text{case} \\ Y_{ij} = 0, j = 1, \dots, m_i & \text{controls} \\ X_{ij1}, \dots, X_{ijp} & \text{primary risk factors} \\ Z_{i1}, \dots, Z_{iq} & \text{matching covariates} \\ & (\text{same value across different people}) \end{array} \right.$$

$$\text{logit } \pi_{ij} = x_i + x_{ij}^T \beta$$

$$P(Y_{i0} = 1, Y_{i1} = 0, \dots, Y_{im_i} = 0) \sum_j Y_{ij} = 1$$

First evaluate the joint probability:

$$P(Y_{i0}=1, Y_{i1}=0, \dots, Y_{im_i}=0)$$

$$= \frac{\exp(\alpha_i + x_{i0}^T \beta)}{1 + \exp(\alpha_i + x_{i0}^T \beta)} \prod_{j=1}^{m_i} \frac{1}{1 + \exp(\alpha_i + x_{ij}^T \beta)}$$

$$= \frac{\exp(\alpha_i + x_{i0}^T \beta)}{\prod_{j=0}^{m_i} \{1 + \exp(\alpha_i + x_{ij}^T \beta)\}}$$

Can derive probabilities for $Y_{i0}=0, Y_{i1}=1, Y_{i2}=0 \dots$
and all other configurations.

$$P\left(\sum_j Y_{ij} = 1\right) = \sum_{j=1}^{m_i} \frac{\exp(\alpha_i + x_{ij}^T \beta)}{\prod_{k=0}^{m_i} \{1 + \exp(\alpha_i + x_{ik}^T \beta)\}}$$

$$P(Y_{i0}=1, Y_{i1}=0, \dots, Y_{im_i}=0 \mid \sum_{j=1}^{m_i} Y_{ij}=1)$$

$$= \frac{\exp(\alpha_i + x_{i0}^T \beta)}{\sum_{j=0}^{m_i} \exp(\alpha_i + x_{ij}^T \beta)} = \frac{\exp(x_{i0}^T \beta)}{\sum_{j=0}^{m_i} \exp(x_{ij}^T \beta)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

- After taking log, can estimate β using Newton-Raphson / Fisher Scoring
- Variance is the inverse of the last iteration of NR
Since α_i disappears, we can't say anything about the population.
- Cannot draw conclusions about potential confounders
- Cannot derive absolute risk
(i.e. prob someone gets lung cancer)
- Can only evaluate relative risk for cancer between smoker and non-smoker, adjusting for confounders.

GLMs

- random
- systematic
- link

difference between
IR and Fisher
scoring

generalized Pearson \rightarrow overdispersion

Pearson log-linear and multinomial (table of counts)



reformulation

$$\frac{D}{np} \sim \chi^2_{n-p} \approx \phi \text{ (not necessarily consistent)}$$

Pearson χ^2 can be used to estimate dispersion param.

↓
this one is consistent

derive algs for QL fitting / $\hat{\beta}$ estimate

- how to draw inference?
- asymptotic variance?

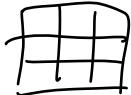
QL is an optimal estimating function within class of
quadratic functions.

model checking

- residual vs. fitted
- score test \approx drop in deviance
asymptotically
- isolated data points (x, y)

influence, Cook's distance

Quasi-symmetric vs. symmetric



$$\log m_{ij} = \lambda_0 + \lambda_i + \lambda_j + \lambda_{ij}$$

$$\log m_{ii} = \lambda_0 + 2\lambda_i$$

↑
no λ_j

Quasi-symmetric: symmetric after getting rid of marginal (row/column) effects