0. Teaching evaluations

1. Interpreting coefs from logistic regression → go back to probabilities
   of Cox & Snell
   also HW 2 Q3
   + deviance test for fit of model

2. Simpson's paradox & reverse — use race example from 199

3. Redelmeier & Singh    +  ~~see~~ HO — Sylvestre & Hanley

4. Survival data   text: 5.4 & 10.7, 8

2.

| | death penalty yes | no | |
|---|---|---|---|
| w | 18 | 141 | 11.88 % |
| ~~victim~~ defendant b | 17 | 149 | 10.24 % |

| white victim d p | | | | black victim d p | | | |
|---|---|---|---|---|---|---|---|
| | yes | no | | | yes | no | |
| w | 19 | 132 | 12% | w | 0 | 9 | 0.2 |
| b | 11 | 52 | 17% | b | 6 | 97 | 5.8 % |

12.6
17.5

1. Binomial regression       $y_i \sim Bin(m_i, p_i(\beta))$    $i = 1, \cdots, n$   ind't

(residual) deviance $D = \sum_{i=1}^{n} \left\{ y_i \log \frac{y_i}{m_i p_i(\hat{\beta})} + (m_i - y_i) \log \left( \frac{m_i - y_i}{m_i (1 - p_i(\hat{\beta}))} \right) \right\}$

is a goodness-of-fit test for the model
$$p_i = p_i(\hat{\beta})$$

(relative to the model $p_i$ unconstrained)

Example HW 2  $\qquad$ logit $p_{ij} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij}$

residual deviance $= 0$  $\qquad$ $p_{ij}(\hat{\beta}) = y_{ij}/m_{ij}$

$$i = 1,2 \; ; \; j = 1,\ldots,6$$

logit $p_{ij} = \mu + \alpha_i + \beta_j$  $\qquad$ no $X^2$

resid deviance $= 20.204$ on 5 df

If $(\alpha\beta)_{ij} \equiv 0$, this $\sim \chi_5^2$

$pr\left(\chi_5^2 > 20.204\right) = 0.001$  "reject $H_0$"

Example 10.15  where resid dev. $= 67.28 \sim \chi_{31}^2$  poor fit

attributed to 2 outliers

(instead of $X^2$)

when those points omitted

$$D = 28.02 \text{ on } 24 \text{ df} \quad \checkmark$$

N.B. This does not work for binary data  Exercise 10.4.1 (a)

$$D = -2 \sum_{i=1}^n \hat{p}_i \log\left(\frac{\hat{p}_i}{p_i}\right) + \log(1-\hat{p}_i)$$

doesn't ~~depend on~~ compare $y_i$ to $\hat{p}_i$

(bec each $y_i = 0$ or 1)

3. Redelmeier & Singh        4650k 'social determinants of health'

- motivation, selection of sample, nominees, controls, winners education
- how do controls work? — same sex, & same age
- response : length of life
- covariates — born in USA, name change, ethnicity, genre, birth yr, sex, age at 1st film, total films, winner/loser or winner/nominee
- time zero
- survivor / that selection bias      lead time bias
- unmeasured confounding — "3 strategies"

Stat. Analysis      primary — KM & log-rank test
regression      Cox PH      note: re quad. & cub. !

Results — baseline characteristics ✓     ⎤
     — career variables ...       ⎥ note: none of these are
     — cause of death         ⎦ 'primary' analyses
\# survival      <u>Figure</u>

- difference in mean life   life expectancy   area under curve   3.9 yrs   $p = 0.003$
and sub-analyses
- 28% reduction in death rate ?? how computed 'Cox' model
- then compared to nominees

Discussion — see marked copy

## 4. Models & methods for survival data

- r.v. $Y$ measures time ($Y \geqslant 0$)

density $f(y)$ on $\mathbb{R}^+$

cdf $F(y) = \Pr(Y \leq y)$

survivor f $= 1 - F(y) = \Pr(Y > y) = S(y)$

<u>hazard</u> f $= h(y) = \dfrac{f(y)}{S(y)}$

cum. haz. $H(y) = \int_0^y h(t)\,dt = -\log S(y)$

i.e. $S(y) = \exp\{-H(y)\}$ $\qquad f(y) = h(y)\exp\{-H(y)\}$

3 examples in §5.4.1.

Weibull f. $S(y) = \exp\left\{-\left(\dfrac{y}{\theta}\right)^\alpha\right\}$ , $\theta, \alpha > 0$

$$f(y) = \frac{\alpha}{\theta}\left(\frac{y}{\theta}\right)^{\alpha-1} \cdot \exp\left\{-\left(\frac{y}{\theta}\right)^\alpha\right\}$$

$$h(y) = \underbrace{\qquad}$$

- censoring : often $\overset{\text{true failure time}}{\underset{\text{not observed}}{\cancel{Y}}}$

random (right) censoring

$$Y_j = \min(Y_j^\circ, C_j) \qquad \begin{array}{l} Y_j^\circ \sim F \\ C_j \sim G \end{array} \text{ indep}^t$$

data $(y_j, \delta_j)$ $j = 1, \dots, n$ $\delta_j = \begin{cases} 1 & obs \\ 0 & censored \end{cases}$

see Figure 3.8

- likelihood f $= \displaystyle\prod_{\delta_j = 1} f(y_j)\{1 - G(y_j)\} \cdot \prod_{\delta_j = 0} S(y_j)\, g(y_j)$

often $f, S$ dep. on $\theta$, but not $G$ $\Rightarrow$ $\ell(\theta) = \displaystyle\sum_{\delta_j = 0} \log f(y_j; \theta)$

Kaplan-Meier estimator of $S(\cdot)$ :

$$\hat{S}(t) = \prod_{i: Y_{(i)} \leq t} \left( \frac{n-i}{n-i+1} \right)^{\delta_{(i)}} , \quad t \leq y_{(n)}$$

$$0 \quad \text{or undefined} \quad t > y_{(n)}$$

$y_{(1)} \leq \dots \leq y_{(n)}$   ordered failure times
$\delta_{(1)}, \dots, \delta_{(n)}$   associated censoring indicators

nb. in Dawson p 197   $\hat{S}(t) = \prod_{j: y_j < t} \left(1 - \frac{1}{r_j}\right)^{d_j}$

$d_j = \delta_{(j)}; r_j = \#$ still alive at time $y_j$ = "risk set for time $y_j$"

can be derived as max lik est., but more obviously is
an extension of eddf $\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^{n} 1\{x_i \leq t\}$

<u>Example</u>   Table 5.3

§10.8.1.   Regression models for censored survival data
$(\underset{\sim}{x}_j, y_j, d_j)$          $f(y_j; x_j, \beta)$ density
    ↑  ↑  ↖ 1, 0
cov.  failure or cens. time      $l(\beta) = \sum_{j=1}^{n} \{d_j \log h(y_j; x_j, \beta) - H(y_j; x_j, \beta)$

PH model   $h(y_j; x_j, \beta) = \underbrace{z(x_j^T \beta)}_{\uparrow} \underbrace{h_0(y_j)}_{\text{baseline hazard}}$
                        increase due
                        to covariates

partial likelihood

$$l_p(\beta) = \prod_{j=1}^{n} \left\{ \frac{\xi(x_j^T\beta)}{\sum_{i \in R_j} \xi(x_i^T\beta)} \right\}^{\delta_j} = \prod_j \left( \frac{\exp{x_j^T\beta}}{\sum_{i \in R_j} e^{x_i^T\beta}} \right)^{\delta_j}$$

usually $\xi(x_j^T\beta) = \exp(x_j^T\beta)$

$$l_p = \sum_{\substack{j=1 \\ \delta_j=1}}^{n} \left[ x_j^T\beta - \log\left\{ \sum_{i \in R_j} \exp(x_i^T\beta) \right\} \right]$$

$$= \sum x_j^T\beta - A_j(\beta)$$

$$l_p' = \sum\left( x_j^T - \frac{B_j}{A_j} \right) \qquad l_p'' = \cdots \qquad (10.62),\ (10.63)$$

Special case: $x_j = \begin{cases} 1 \\ 0 \end{cases}$ group $\begin{matrix} A \\ B \end{matrix}$

$p=1$

$$l_p'(\beta) = \sum_{\delta_j=1} \left\{ x_j - \frac{\sum_{i \in R_j} x_i e^{x_j\beta}}{\sum_{i \in R_j} x_j e^{x_j\beta}} \right\}$$

$$l_p'(\beta) \Big|_{\beta=0} = \sum_{\delta_j=1} \left( x_j - \frac{\sum_{i \in R_j} x_j}{\sum_{i \in R_j} 1} \right) = \sum_{\delta_j=1} \left( x_j - \frac{m_{1j}}{m_{0j}+m_{1j}} \right)$$