

- Normal approximation might be too much a stretch

A suit fitting better for binomial response

- when n_i 's not large enough or $p_x \approx 1$ or 0 → Recall. LM data with response having upper/lower bound, & many observations of the response are close/equal to the bounds.
- Some of these problems could be corrected by using transformation and weighting

eg. $\log(\hat{p}_i/(1-\hat{p}_i)) = X\beta + \varepsilon_i$
 $0 \leq \text{predicted prob.} \leq 1$

for non-constant variance
 for non-constant non-linearity

2/26

Generalized Linear Model for Binomial Data

- Recall: linear model

binomial response $\sim B(n_x, p_x)$
 $\{0, 1, \dots, n_x\} \ni y_x = X\beta + \varepsilon$

equivalent in LM

model description 1: $Y = X\beta + \varepsilon$, $\varepsilon \sim N(0, \sigma^2 I)$

model description 2: $Y \sim N(X\beta, \sigma^2 I)$

response

隨機

$Y \sim N(\mu_x, \Sigma_x)$, $\mu_x \in (-\infty, \infty)$

covariates

規律

coefficients (parameters) known functions (base functions)

$X\beta = \sum_{i=1}^p \beta_i \cdot h_i(X_1, \dots, X_m) \equiv \eta_x(\beta)$

build the link b/w parameters in Y & $X\beta$

$\mu_x = \eta_x$, $\Sigma_x = \sigma^2 I$

matrix form

functional form

expressed as a function or vector

can be assigned a function of x (with parameters) if necessary

a function of x with a linear structure b/w β_i 's & β_i 's

Q: which description can be generalized to binomial data?

of main interest in regression

- 3 components in a generalized linear model (binomial example)

response

隨機

$y_x \sim B(n_x, p_x)$

covariates

規律

$X\beta = \sum_{i=1}^p \beta_i \cdot h_i(X_1, \dots, X_m) \equiv \eta_x(\beta)$

$\mu_x = E(y_x) = n_x p_x$
 $\sigma_x^2 = \text{Var}(y_x) = n_x p_x (1 - p_x)$
 $= [\mu_x(n_x - \mu_x)]/n_x$

σ_x^2 is a function of μ_x (cf. LM)

$\mu_i = n_i p_i$
 $p_i = \mu_i / n_i$
 $\mu_i = n_i p_i$

Build the link b/w the parameters in y_x & $X\beta$

link function g : monotone and differentiable such that $\eta_x = g(p_x) \Rightarrow p_x = g^{-1}(\eta_x)$ [for binomial, $g: (0, 1) \rightarrow (-\infty, \infty)$]

- Common choices of link function for binomial data

canonical link

Logit: $\eta_x = \log(p_x / (1 - p_x))$

logistic regression $p_x \xrightarrow{\text{logit}} \eta_x$
 $\eta_x \xrightarrow{\text{logistic}} p_x$

Probit: $\eta_x = \Phi^{-1}(p_x)$, where Φ is the cdf of $N(0, 1)$

the pdfs of Φ & logistic are symmetric about 0

Complementary log-log: $\eta_x = \log(-\log(1 - p_x))$

Note. g^{-1} corresponds to a cdf with heavy tail

Note.

Logit is close to the complementary log-log when p_x is small

Logit is close to probit when $0.1 < p_x < 0.9$

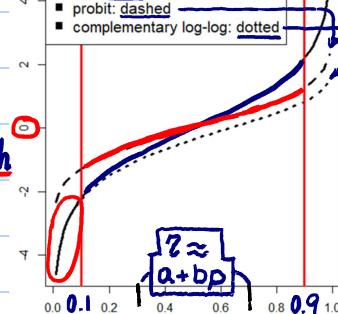
(exercise: compare the 3 functions)

a linear transformation

a non-linear function \Rightarrow relationship b/w μ & x is not linear any more

different choice of g results in different interpretation of β

logit: solid
 probit: dashed
 complementary log-log: dotted



- widely
appli-
cable

textbook,
Appendix A

- good statistics
- asymptotic properties

- under
logit
link

$$\eta = g(p) = \log\left(\frac{p}{1-p}\right)$$

• $p = g^{-1}(z)$
logistic = $\frac{e^z}{1+e^z}$

under
any
links

$$\bullet \underline{1-p} = \frac{1}{1+e^2}$$

Y_i indep

$$\underline{p} : (1-p) = \underline{e^2} : \underline{e^4}$$

- LNp. 2-17, pmf of $B(n_i, p_i)$
 \therefore independent

↳ $(X^T X)^{-1} X^T Y$ (good projection properties)

Why? check
model des-
cription 1
in LNp 3-4

- only needs likelihood

(But, no cription 1
in / N-3-11

partial derivative
of each θ_r ,
 $r = 1, \dots, p$

project-
tion explanation)

$$r=1, \dots, p$$


for one
covariate
class
 $(\underline{x}_i, \underline{y}_i)$

— canonica

21

check *

in $LN_{p.3-5}$

 $\partial i / \partial p_i \left[-\sigma^{-1} \right]$

for any links

$$= \frac{y_i}{p_i} - \frac{n_i - y_i}{1 - p_i} = \frac{y_i - n p_i}{p_i (1 - p_i)}$$

$$\frac{\partial p_i}{\partial \eta_i} = \frac{\partial}{\partial \eta_i} \frac{e^{\eta_i}}{1+e^{\eta_i}} = \frac{e^{\eta_i}(1+e^{\eta_i}) - e^{\eta_i} \cdot e^{\eta_i}}{(1+e^{\eta_i})^2}$$

- Usually no explicit formula for M

In LM,

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

In LM, for non-constant variances modeled by unknown parameters check * in LNp.3-4)

Why? Note that $\text{Var}(y_x)$ is a function of $E(y_x)$ (Lnp. 3-5)



made by S.-W. Cheng (NTHU, Taiwan)

Test and Confidence Interval ← Inference

- RSS (in LM) and deviance (in GLM)

residual
 $y_i - \hat{y}_i$

➤ Recall: in linear model, RSS play a critical role in inference

➤ Q: How to evaluate whether y_i & \hat{y}_i are close enough in GLM?

➤ Q: In GLM, what is the concept similar to RSS?

check
LNp.1-13
Test
in LM

➤ consider two models L and S l, s : # of parameters in $X\beta$

- a larger model L : l parameters and likelihood \mathcal{L}_L
- a smaller model S : s ($s < l$) parameters and likelihood \mathcal{L}_S
- S is nested in L ($S \subset L$), e.g., any joint distribution of data in $S \in L$

likelihood ratio
 $-2 \log \left(\frac{\mathcal{L}(\hat{\beta}_S)}{\mathcal{L}(\hat{\beta}_L)} \right)$

$\dim = s$
 $\text{span}(X_S) \subset \text{span}(X_L)$
 $\dim = l$

$L: X_L \beta_L = \eta_L = g(\eta_L) \leftarrow H_0: H_1$
 \downarrow
 $\& H_0: A\beta_L = c$
 $S: X_S \beta_S = \eta_S = g(\eta_S) \leftarrow H_0$
 $A\beta_L \neq c$

➤ To test $H_0: S$ (say, $A\beta = c$) vs. $H_1: L \setminus S$, likelihood methods suggests the likelihood ratio statistics:

$$2[l(\hat{\beta}_L) - l(\hat{\beta}_S)] = 2 \log \left[\frac{\mathcal{L}(\hat{\eta}_L = X_L \hat{\beta}_L)}{\mathcal{L}(\hat{\eta}_S = X_S \hat{\beta}_S)} \right] \stackrel{a}{\sim} \chi_{df_{L \setminus S}}^2$$

log-likelihood (LNp.3-6) where $df_{L \setminus S} = \dim(L) - \dim(S) = l - s$. MLE under L MLE under S approximately or asymptotically (n_i 's large)

■ Suppose that the larger model L is saturated.

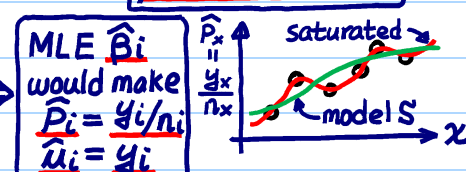
- For R covariate classes with distinct x_i 's, a saturated L has R parameters in β .

• Say,

$$X = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}_{R \times R}, \beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \Rightarrow u_i = n_i p_i = n_i g'(\eta_i) = n_i g'(\beta_i), i=1, \dots, R$$

can freely change in estimation $(X\beta)_i$

an identifiable model with maximum # of parameters



For other saturated X^* , $\text{span}(X^*) = \text{span}(X)$

We have $\hat{y}_{i,L} = y_i$ under L , and the LR test statistic becomes:

$$D_S = 2 \sum_{i=1}^k \left\{ y_i \log(y_i / \hat{y}_{i,S}) + (n_i - y_i) \log[(n_i - y_i) / (n_i - \hat{y}_{i,S})] \right\}$$

$2[l(\hat{\beta}_L) - l(\hat{\beta}_S)]$ # of success $y_i \leftrightarrow \hat{y}_{i,S}$ ≈ 1 , if $y_i \approx \hat{y}_{i,S}$ # of failure $n_i - y_i \leftrightarrow n_i - \hat{y}_{i,S}$ ≈ 1 , if $y_i \approx \hat{y}_{i,S}$

small D_S better fit $\hat{y}_i \leftrightarrow y_i$

➤ D_S is called deviance of S , which plays a role similar to RSS

■ Since the saturated model fits as well as any model can fit, the deviance D_S measures how close the (smaller) model S comes to the perfection (i.e., $D=0$ under saturated model).

can do test !!!

Deviance can be treated as a measure of goodness-of-fit In LM, can use (i) $\hat{\sigma}^2 = \text{RSS}/n-p$ (ii) $R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$

Compare S with saturated model \leftarrow most complicated \leftrightarrow compare S with intercept-only model \rightarrow simplest

prefer S with small D_S & fewer parameters