

Logistic Regression I

Summary of the last lecture

- Types of residuals
- Residual plots
- Over-dispersion

Key terms of this lecture

- Logistic regression
 - Terminology
 - Ungrouped data/Grouped data

Reading

- McCullagh and Nelder (1989) Chapter 7
- Dobson and Barnett (2008) Chapter 4

1

Binomial Distribution

- Example: Budworm data. This example comes from Venables and Ripley's *Modern Applied Statistics with S*, Springer 4th edition, 2002. [See the attached pages for detail.]
 - Batches of 20 moths subjected to increasing doses of a poison, "success" = death.
 - Data is grouped: for each of 6 doses (1.0, 2.0, 4.0, 8.0, 16.0, 32.0 mg) and each of male and female, we have 20 moths. *in each group.*
 - There are 12 covariate patterns. *or there are $6 \times 2 = 12$ groups based on different values of X 's.*

dose	1	2	4	8	16	32
Male	1	4	9	13	18	20
Female	0	2	6	10	12	16

The number of moths that died in each group is shown in the table. Data suggests that the probability of dying increases with dose.

2

Terminology

- Odds for a single moth to die

$$\text{odds}(p) = \frac{p}{1-p}$$

- Log-odds (a.k.a logit)

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right)$$

3

Terminology (cont'd)

- Between two groups, e.g., exposed vs. unexposed, Trt1 vs. Trt2)
 - In budworm example, consider the male budworms receiving doses 1 and 2 at which $y_1 = 1$ and $y_2 = 4$ moths died out of $m = 20$.
 - Risk Difference: $RD = p_1 - p_2$ (e.g., $0.05 - 0.2 = -0.15$) $\frac{1}{20} - \frac{4}{20} = 0.05 - 0.2$
 - Relative Risk: $RR = p_1/p_2$ (e.g., $0.05/0.2 = 0.25$)
 - Odds of dying: $\text{odds}(p_1) = p_1/(1-p_1)$ and $\text{odds}(p_2) = p_2/(1-p_2)$
 - Odds ratio: $OR = \text{odds}(p_1)/\text{odds}(p_2)$ (e.g., $(0.05/0.95)/(0.2/0.8) = 0.21$) *odds of dying at dose 1 related to dose 2*
 - Log-Odds ratio: $\text{Log-OR} = \log(OR)$ (e.g., $\log(0.21) = -1.56$)

4

Terminology (cont'd)

- Interpretation:
 - RD : hard to explain a difference of -0.15 in probability
 - RR and OR : A ratio of 0.25 in probability is easier to interpret.
 - * $RR = 0.25$: The risk of dying from dose 1 is only 1/4 of the risk of dying from dose 2.
 - * Thus, having $RR < 1$ means that dose 1 is less poisonous than dose 2.
 - * Note: 0.25 does not say anything about the risk of dying from dose 2.
 - Log- OR $OR \in (0, \infty)$; $\text{Log-}OR \in (-\infty, \infty)$. Useful in interpretation of parameters in logistic regression.

5

Grouped vs. Ungrouped Data

- Grouped data (budworm example, seed example)
 - Data are presented by distinct covariate values (e.g., dose-gender seed-extract combination)
 - After counting number of responses in each category, we have a binomial r.v. (i.e. Response is the number of success)
- Ungrouped data
 - Each record represents one individual. Response is binary (0 or 1)
- Example. Seed data
 - Both results are exactly the same. *Grouped data vs. ungrouped data. Estimates and standard errors are the same, but the measures of goodness of fit differ. See example: Senility and Waive page 140, 3rd ed.*
 - Differences:
 - * DFs are different, so is GOF
 - * Overdispersion be found in grouped data. *Further grouping: the Hosmer-Lemeshow statistic χ^2_{HL} . By grouping values of $\hat{\pi}_i$, so that the total number of observations per group are approximately equal*

Budworm Example

- Goal: describe the dependence of the mortality of budworms on sex S and the applied dose D .
 - i.e. explain $E[Y]$ as a function of the covariates.
- Data:
 - D : logarithm of dose (1,2,4,8,16,32 (μg)) of cypermethrin (continuous covariates)
 - S : sex (male, female) (class or factor covariates)
- Distribution assumption: $y_{ij} \sim \text{Binomial}(n_{ij}, p_{ij})$, for $i = 1, 2, j = 1, \dots, 6$, with $E[Y_{ij}] = n_{ij}p_{ij}$ where $n_{ij} = 20$ known in this example.
- Thus, the question is the relationship between p_{ij} and covariates.

7**Budworm Example (cont'd)**

- Exploratory Plots: see R and SAS code, budworm.R and budworm.SAS
 - Observed proportions

$$\hat{p}_{ij} = \frac{y_{ij}}{n_{ij}}.$$

- Empirical logits

$$\text{elogeit}_{ij} = \log \left(\frac{y_{ij} + 0.5}{n_{ij} - y_{ij} + 0.5} \right)$$

- What do plots suggest? [See the attached plots.]

8

Budworm Example (cont'd)

- Components of GLM logistic regression

- Random component

$$Y_{ij} \sim \text{Bin}(n_{ij}, p_{ij}), \quad n_{ij} = 20$$

$$\text{or } Y_i \sim \text{Bin}(n_i, p_i), \quad n_i = 20$$

- Linear predictor (Systematic component)

$$\eta_i = \beta_0 + \beta_1 \text{Sex} + \beta_2 \log \text{dose}$$

- Link function

$$\eta_i = \log \frac{p_i}{1-p_i} \quad \text{or} \quad \log \frac{p_i}{1-p_i} = \eta_i$$

- Thus, the model would be,

$$\log \frac{p_i}{1-p_i} = \beta_0 + \beta_1 \text{Sex} + \log \text{dose}$$

9

Budworm Example (cont'd)

- Parameter estimation

- MLE of β using IWLS.

- Interpretation of Parameters

- β_0 : log-odds when $\text{Sex} = 0$ (Female) and $\log \text{dose} = 0$ (dose = 1)
- β_1 : log-OR when Sex is changed from Female (=0) to male (=1), holding logdose constant
- $\beta_0 + \beta_1$: log-odds when $\text{Sex} = 1$ (male) and $\log \text{dose} = 0$ (dose = 1)
- β_2 : log-OR when logdose is increased by one unit, while holding Sex unchanged.

10

- Revisit:
 - Odds

$$\text{odds}(p_{ij}) = \frac{p_{ij}}{1 - p_{ij}} = \exp(\text{logit}(p_{ij})) = \exp(\eta_{ij}),$$

Thus

$$p_{ij} = \frac{\exp(\text{logit}(p_{ij}))}{1 + \exp(\text{logit}(p_{ij}))} \\ = \frac{\exp(\eta_{ij})}{1 + \exp(\eta_{ij})}$$

the fitted probabilities: $\hat{p}_{ij} = \frac{\exp(\hat{\eta}_{ij})}{1 + \exp(\hat{\eta}_{ij})}$, $\hat{\eta}_{ij} = \mathbf{x}_{ij}^T \hat{\beta}$

$\hat{\eta}_{ij}$ vs log dose — linear

\hat{p}_{ij} vs log dose — non-linear.

11

Budworm Example (cont'd): Interpretation

- Interpretation of a continuous covariate
 - Using GLM for Binomial distribution, let $D = \log(\text{dose})$,

$$\text{logit}(p_{i,d+1}) = \beta_0 + \beta_1 S + \beta_2 (D + 1)$$

$$\text{logit}(p_{i,d}) = \beta_0 + \beta_1 S + \beta_2 D$$

thus, β_2 is a change in logits per unit increase of log-dose. i.e.

$$\text{odds}(p_{i,D+1}) = \exp(\text{logit}(p_{i,D+1})) \\ = \exp(\beta_0 + \beta_1 S + \beta_2 (D+1))$$

thus, in example, increasing the log-dose by one unit increases the odds of death by a factor of $\exp(\beta_2)$

$$\begin{aligned} \text{e.g., Wald CI for } \beta_2: \hat{\beta}_2 \pm 1.96 \text{SE}(\hat{\beta}_2) \\ \text{e.g., Wald CI for OR } \exp(\text{CI for } \beta_2) \end{aligned} \quad \begin{aligned} &= \exp(\beta_2) \exp(\beta_0 + \beta_1 S + \beta_2 D) \\ &= \exp(\beta_2) \text{odds}(p_{i,D}) \\ &\Rightarrow \exp(\beta_2) = \frac{\text{odds}(p_{i,D+1})}{\text{odds}(p_{i,D})} = \text{OR} \end{aligned}$$

e.g., Wald CI for OR $\exp(\text{CI for } \beta_2)$

$$= \exp(\hat{\beta}_2 \pm 1.96 \text{SE}(\hat{\beta}_2))$$

$$\beta_2 = \log \text{OR}$$

12

This CI Excluding unity (1) indicates β_2 is significantly different from zero.

Budworm Example (cont'd): Interpretation

- Interpretation of a factor

- Same as one for a continuous covariate

$$\begin{aligned} \overset{\text{Male}}{\uparrow} \text{logit}(p_{M,D}) - \overset{\text{Female}}{\uparrow} \text{logit}(p_{F,D}) &= (\beta_0 + \beta_1(1) + \beta_2 D) - (\beta_0 + \beta_1(0) + \beta_2 D) \\ &= \beta_1 \end{aligned}$$

thus, the *OR* of dying as a male moth relative to a female moth is $\exp(\beta_1)$

i.e. the odds of dying for males is about $\exp(\beta_1)$ times the odds for females.

- Confidence interval for β_1 e.g., Wald CI: $\hat{\beta}_1 \pm 1.96 SE(\hat{\beta}_1)$
- Confidence interval for the odds ratio (*OR*) is.. $\exp(\hat{\beta}_1 \pm 1.96 SE(\hat{\beta}_1))$

13

Budworm Example (cont'd): Probability

- Estimation (or prediction) of the probability of dying at certain levels of the covariates

$$\text{logit}(p_{ij}) = \eta = \beta_0 + \beta_1 S + \beta_2 D.$$

- (1) Estimate of log-odds at female moth to die at dose 27 (i.e. $\log(27) = 3.3$).

$$\hat{\eta} = \hat{\beta}_0 + \hat{\beta}_1(S = 0) + \hat{\beta}_2(D = 3.3).$$

- Confidence interval for η (log-odds, here)

$$\hat{\eta} = \mathbf{C}\hat{\beta},$$

where $\mathbf{C} = (1, S, D) = (1, 0, 3.3)$. Using Wald Statistic,

$$(\mathbf{C}\hat{\beta} - \mathbf{C}\beta)^T \{\hat{V}(\mathbf{C}\hat{\beta})\}^{-1} (\mathbf{C}\hat{\beta} - \mathbf{C}\beta) \sim \chi^2_1$$

14

Budworm Example (cont'd): Model Comparison

- Consider an interaction model

(1) Test $H_0: \beta_0 = \beta_1 = \beta_2 = \beta_3$ *Testing all parameters are equal*

$$\Leftrightarrow C\beta = 0, \quad C = \begin{pmatrix} 1 & -1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 1 & 0 & 0 & -1 \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}, \quad \hat{V}(\hat{\beta}) = \begin{pmatrix} 0.30548 & * & * & * \\ & 0.60577 & & \\ & & 0.05812 & \\ * & & & 0.15772 \end{pmatrix}$$

Then $(C\hat{\beta} - C\beta)^T \hat{V}(\hat{\beta})^{-1} (C\hat{\beta} - C\beta) = (C\hat{\beta} - 0)^T \hat{V}(\hat{\beta})^{-1} (C\hat{\beta} - 0) \sim \chi^2(3)$

(2) Test $H_0: \beta_3 = 0$.

To compare the main effects model and the interaction model and test importance of covariates.

Wald test: $\hat{\beta}_3 = 0.5091$, $SE(\hat{\beta}_3) = 0.3895$, 95% CI: $(-0.2543, 1.2726)$

$H_0: \beta_3 = 0$, $\chi^2_{calc(1)} = 1.71$, $p\text{-value} = 0.1912$, non-significant

Deviance: $D_2 - D_1 = 6.7571 - 4.9937 = 1.7634 \approx LR$ (In fact, they are equivalent).

Pearson χ^2 : $\chi^2_{calc(1)} = 5.3060 - 3.5047 = 1.8013$, non-sig.

$LR = 2 [(-105.7388) - (-106.6204)] = 2 (0.8816) = 1.7632$ non-significant

Here, when $q=1$, $D^*(Y, \hat{\mu}) = D(Y, \hat{\mu})/q \approx D(Y, \hat{\mu})$, so $LR = Dev.$

Budworm Example (cont'd): GOF

- Using Deviance, For model with interaction

$$D(Y, \hat{\mu}) = 4.9937, \quad df = 12 - 4 = 8$$

Compare this model with the saturated model with the n parameters, the test is not significant, indicating a good fit

- Using Pearson's χ^2

$$\chi^2 = \sum_{i=1}^n \frac{(\tilde{Y}_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)} = 3.5047, \quad df = 8$$

where $\tilde{Y}_i = Y_i/m$ is the response variable for Binomial proportion. *Conclusion is the same*

For binomial distribution, $V(\hat{\mu}_i) = \hat{\mu}_i(1 - \hat{\mu}_i)/m_i$

$$\chi^2 = \sum_i \frac{(Y_i - m_i \hat{\mu}_i)^2}{m_i(1 - \hat{\mu}_i)\hat{\mu}_i}$$

Thus, the CI for η is

$$\hat{\eta} - \sqrt{\chi_{0.05,1}^2 \hat{V}(C\hat{\beta})} \leq \eta \leq \hat{\eta} + \sqrt{\chi_{0.05,1}^2 \hat{V}(C\hat{\beta})}.$$

(2) Estimation of p

$$\hat{p} = \frac{\exp(\hat{\eta})}{1 + \exp(\hat{\eta})}$$

and the CI for p is

$\because \eta = C\beta$, $\text{Var}(\hat{\eta}) = \hat{V}(C\hat{\beta}) = C \hat{V}(\hat{\beta}) C^T$. $\hat{V}(\hat{\beta})$ can be obtained from SAS budworm.sas, by model ndead/ntotal=Sex logdose / dist=binomial covb;
i.e., $\hat{V}(\hat{\beta}) = \begin{pmatrix} 0.21951 & -0.09875 & -0.07575 \\ -0.09875 & 0.12601 & 0.01875 \\ -0.07575 & 0.01875 & 0.03576 \end{pmatrix}$.

\therefore 95% CI for $\eta = C\beta$ is $\hat{\eta} \pm \sqrt{\chi_{0.05,1}^2 \hat{V}(C\hat{\beta})} = \hat{\eta} \pm 1.96 \sqrt{C \hat{V}(\hat{\beta}) C^T}$
and 95% CI for p is $\frac{\exp(\hat{\eta} - 1.96 \sqrt{C \hat{V}(\hat{\beta}) C^T})}{1 + \exp(\hat{\eta} - 1.96 \sqrt{C \hat{V}(\hat{\beta}) C^T})} < p < \frac{\exp(\hat{\eta} + 1.96 \sqrt{C \hat{V}(\hat{\beta}) C^T})}{1 + \exp(\hat{\eta} + 1.96 \sqrt{C \hat{V}(\hat{\beta}) C^T})}$

15

Budworm Example (cont'd): Interaction

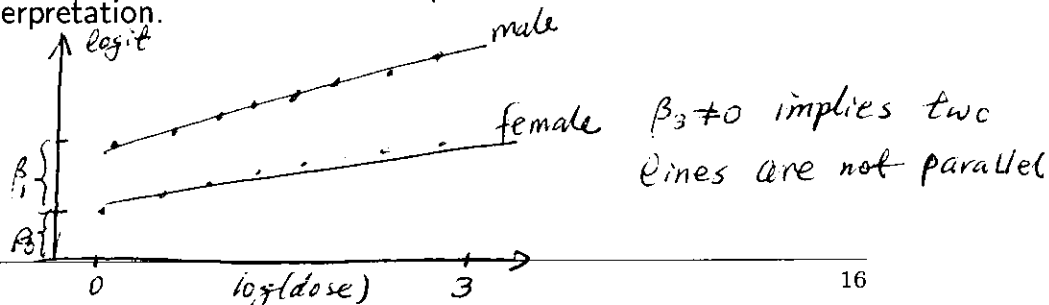
- Consider a model with interaction

$$\text{logit}(p) = \eta = \beta_0 + \beta_1 S + \beta_2 D + \beta_3 S \times D.$$

- Interpretations of parameters.

- β_0 : log-odds at $S=0$ (Female) and $D=0$ (dose=1)
- β_1 : log-OR of dying as a male ($S=1$) moth relative to a female.
- β_3 : Interaction effect, log-OR of dying as a male moth relative to a female moth at a fixed D level when $D=1$

- Graphical interpretation.



16

