

# STA 303/1002-Methods of Data Analysis II

Sections L0101& L0201, Winter 2019

**Dr. Shivon Sue-Chee**



February 4-8, 2019

## STA 303/1002: Week 5-Generalized Linear Models

- ▶ Case Study III: The Donner Party Example
- ▶ Generalized Linear Models
  - ▶ What is a Generalized Linear Model?
  - ▶ Common link functions
  - ▶ What is a Binary Logistic Regression Model?
  - ▶ Maximum likelihood estimation of  $\beta$ 's
- ▶ Case Study III Example
  - ▶ Data and Questions
  - ▶ Estimated Model
  - ▶ Interpretations

## Case Study III: The Donner Party Example

- ▶ Background: (D.K. Grayson, Journal of Anthropological Research, 1990: 223-42 )
  - ▶ In mid 19th century, a group of 86 American pioneers headed out from Missouri toward California in a wagon train.
  - ▶ Due to a combination of harsh weather, unsuitable travel equipment and divisions with the party, the group got stuck in the Sierra Nevada mountain range.
  - ▶ They had planned to arrive safe and sound in September but those who survived did not make it there until the following March.

▶ Question: Who survived?- Men? -Older pioneers?

▶ Data:

- ▶ age
- ▶ sex
- ▶ outcome: survived or not

▶ AIM: Study the odds of survival

## Case Study III: Model

- ▶ Response:  $Y_i$ - a binary variable (eg., survived or died)
- ▶ Predictor:  $X_i$ - eg., age, sex of  $i$ th pioneer
- ▶ Model: BINARY LOGISTIC REGRESSION

$$Y_i|X_i = \begin{cases} 1 & \text{if response is in category of interest} \\ 0 & \text{otherwise} \end{cases}$$

$$Y_i|X_i \sim \text{Bernoulli}(\pi_i)$$

$$0 \leq \pi \leq 1$$

Then:

- ▶  $E[Y_i|X_i] = \pi_i$  and  $\text{Var}(Y_i|X_i) = \pi_i(1 - \pi_i)$
- ▶ A logistic regression model is an example of a **Generalized Linear Model**.

## Generalized Linear Models

- ▶ Have: · response,  $Y$  and  
· a set of explanatory variables  $X_1, \dots, X_p$
- ▶ Want: Model  $E(Y)$  as a linear function in the parameters, ie.,

$$g(E(Y)) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p = \mathbf{X}\beta$$

- ▶ Key idea: Choice of the link function,  $g$  such that

$$g(E(Y)) = \mathbf{X}\beta$$

transformation  
 $\downarrow g(y)$

$$E[g(y)] \neq g[E(y)]$$

## Some Link Functions

Let  $E(Y) = \mu$ .

Link	Function	Usual distribution of $Y X$
Identity	$g(\mu) = \mu$	Normal — General Linear Model
Log	$g(\mu) = \log \mu, \mu > 0$	Poisson (count data)
Logit	$g(\mu) = \log \left( \frac{\mu}{1-\mu} \right), 0 < \mu < 1$	Bernoulli (binary), Binomial

Note: Link function,  $g(\cdot)$  is a function of  $\mu = E(Y)$ , the mean of  $Y$ , and not a transformation of the data.

$$\log \left( \frac{\pi}{1-\pi} \right)$$



## GLMs vs Transforming the data

- ▶ Transform Y so it has an approximate normal distribution with constant variance. Common variance stabilizing transformations (Weisberg, 3rd ed, p. 179):
  - ▶  $\sqrt{Y}$ : mild transformation; used when  $\text{Var}(Y|X) \propto E(Y|X)$  as for Poisson data
  - ▶  $\log(Y)$ : most common; if  $\text{Var}(Y|X) \propto [E(Y|X)]^2$  or errors behave like percentage of Y.
  - ▶  $1/Y$ : used when responses are mostly close to 0, but some large values occur.
- ▶ As GLM (Agresti, p. 117):
  - ▶ distribution of Y not restricted to Normal
  - ▶ model parameters describe  $g[E(Y)]$  rather than  $E(g(Y))$  as in transformed data approach
  - ▶ GLMs provide a unified theory of modelling that encompasses the most important models for continuous and discrete variables.

Box-Cox

## LOG ODDS, ODDS, ODDS RATIO

- ▶ Let  $\pi = P(\text{"success"})$ ,  $0 < \pi < 1$ .

- ▶ The ODDS in favour of "success" is:

$$(0, \infty)$$

$$\frac{\pi}{1 - \pi} = \frac{P(\text{"success"})}{P(\text{"failure"})}$$

- ▶ Then the LOG ODDS is:

$$(-\infty, \infty)$$

$$\log \left( \frac{\pi}{1 - \pi} \right)$$

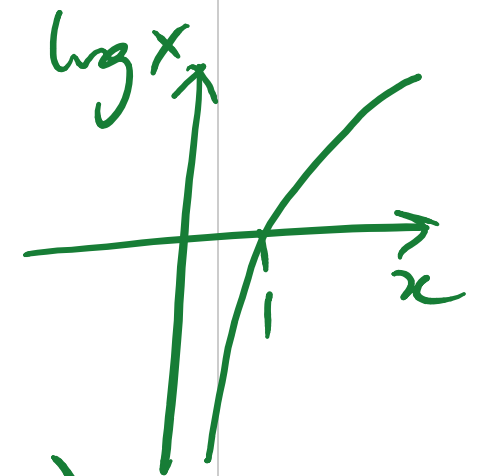
- ▶ An ODDS RATIO is a ratio of ODDS.

$$(0, \infty)$$

$$\frac{\pi_1 / (1 - \pi_1)}{\pi_2 / (1 - \pi_2)} = \frac{\pi_1}{\pi_2} \left( \frac{1 - \pi_2}{1 - \pi_1} \right)$$

$$\frac{1}{3}$$

$$\frac{5}{1}$$





## Binary Logistic Regression

**AIM** ▶  $E(Y|X) = \pi$

▶  $\text{Var}(Y|X) = \pi(1 - \pi)$ . Notice that variance is not constant!

▶ Logistic regression model:

logit  $\left[ \log \left( \frac{\pi}{1 - \pi} \right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p \right] \quad (1)$

*mean kernel function  
linear in the parameters*

▶ Linear predictor:

$$\eta = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

▶ **LOGISTIC FUNCTION**: Find by inverting equation (1)

$$\pi(\eta) = \frac{e^\eta}{1 + e^\eta} = \frac{1}{1 + e^{-\eta}}$$

$$= \frac{e^\eta / e^\eta}{1/e^\eta + e^\eta / e^\eta}$$

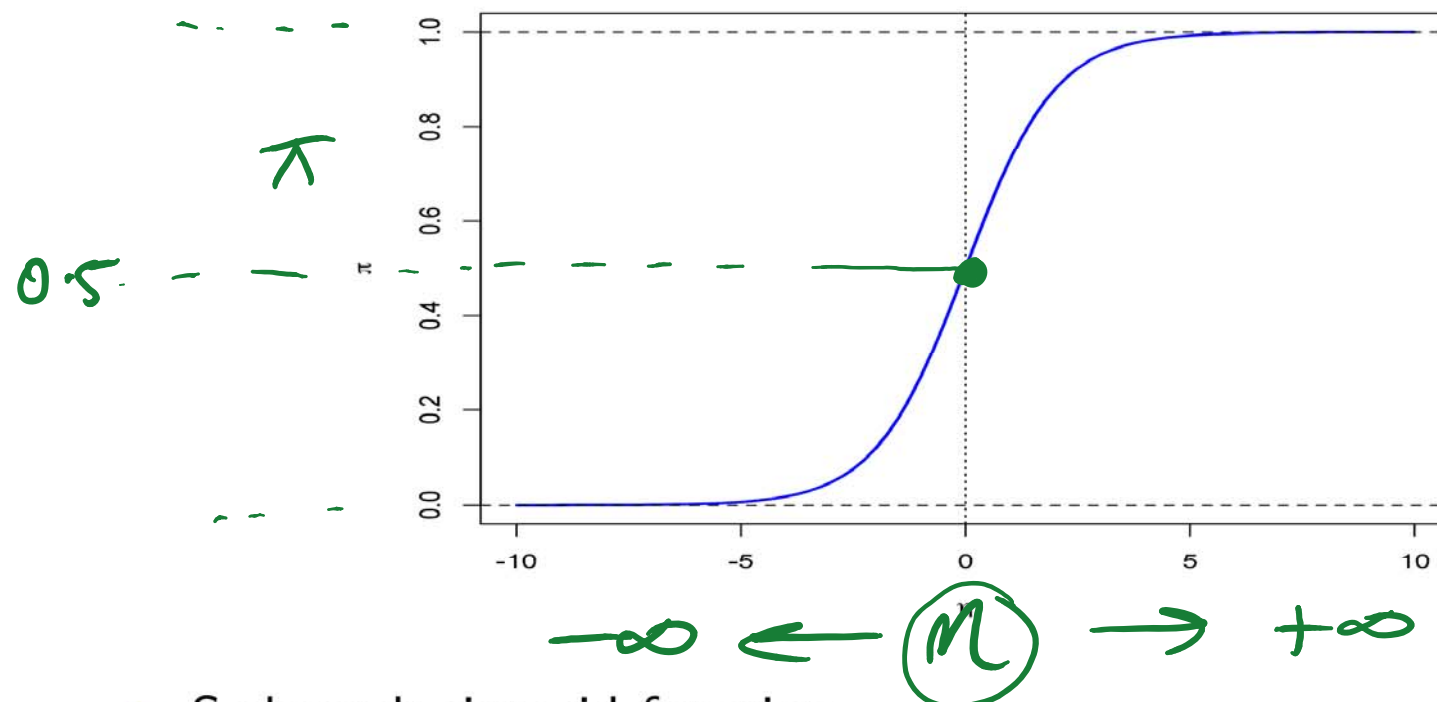
$$\log \left( \frac{\pi}{1 - \pi} \right) = \eta$$

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

$$\pi = e^\eta - \pi e^\eta$$

What does the logistic function look like?

► LOGISTIC FUNCTION:  $\pi = \frac{e^\eta}{1+e^\eta} = \frac{1}{1+e^{-\mathcal{M}}}$



$$\begin{array}{l} \mathcal{M} \rightarrow +\infty \\ \pi \rightarrow 1 \\ \hline \mathcal{M} \rightarrow -\infty \\ \pi \rightarrow 0 \\ \hline 0 < \pi < 1 \end{array}$$

- S-shaped; sigmoid function
- Horizontal asymptotes at 0 and 1; the logistic function,  $\pi(\eta)$  varies between 0 and 1

## Binary Logistic Regression Model

loge  
ln

$$\log \left( \frac{\pi_i}{1 - \pi_i} \right) = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip}, \quad i = 1, \dots, n$$

- ▶ Log-odds,  $\log(\pi/(1 - \pi))$  are between  $-\infty$  and  $\infty$  (good characteristic of a link function)
- ▶ As  $\pi_i$  (the probability of “success”) increases, odds of success and log-odds increase
- ▶ Predicts the natural log of the odds for a subject being in one category or another
- ▶ Regression coefficients can be used to estimate odds ratio for each of the independent variables
- ▶ Tells which predictors can be used to determine if a subject was in a category of interest

# How to estimate the parameter coefficients?

## Maximum Likelihood Estimation

- ▶ Data:  $Y_i = \begin{cases} 1 & \text{if response is in category of interest} \\ 0 & \text{otherwise} \end{cases}$
- ▶ Model:  $P(Y_i = y_i) = \pi_i^{y_i} (1 - \pi_i)^{1-y_i}$
- ▶ Assume: The  $n$  observations are independent
- ▶ Joint density:

$$P(Y_1 = y_1, \dots, Y_n = y_n) = \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i}$$

where

$$\pi_i = \frac{\exp(\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip})}{1 + \exp(\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip})} = \frac{e^M}{1 + e^M}$$

$$\text{and } 1 - \pi_i = \frac{1 - \frac{e^M}{1 + e^M}}{1 + e^M} = \frac{1 + e^M - e^M}{1 + e^M} = \frac{1}{1 + e^M} = (1 + e^M)^{-1}$$

## Maximum Likelihood Estimation

$y_i$

$$\log ab = \log a + \log b$$

- **Likelihood function:** Plug in observed data and think of the joint density as a function of  $\beta$ 's-

$$\mathcal{L}(\beta_0, \dots, \beta_p) = \prod_{i=1}^n \pi_i(\beta)^{y_i} (1 - \pi_i(\beta))^{1-y_i}$$

- **Log-likelihood function:**  $= \prod_{i=1}^n \left( \frac{e^{\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}}}{1 + e^{\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}}} \right)^{y_i} \left[ \left( 1 + e^{\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}} \right)^{-1} \right]^{1-y_i}$

$$\begin{aligned} \log \mathcal{L}(\beta_0, \dots, \beta_p) &= \sum_{i=1}^n [y_i(\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}) \\ &\quad - y_i \log(1 + \exp(\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip})) \\ &\quad - (1 - y_i) \log(1 + \exp(\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}))] \end{aligned}$$

- Maximize the log-likelihood:

$$(\hat{\beta}_0, \dots, \hat{\beta}_p) = \arg \max \{ \log \mathcal{L}(\beta_0, \dots, \beta_p) \}$$

## MLE solution methods

- ▶ No explicit expression exists for the maximum likelihood estimators  $-(\hat{\beta}_0, \dots, \hat{\beta}_p)$ .
- ▶ Two iterative numerical solution methods are:
  - (1) Newton-Raphson algorithm
  - (2) Fisher scoring or Iteratively Re-weighted Least Squares (IWLS). This is done in `glm()`.

## Large-sample properties of MLEs

If model is correct, and sample size is large enough, as  $n \rightarrow \infty$

1. MLEs are <sup>nearby</sup> unbiased
2. MLEs have minimum variance
3. MLEs are Normally distributed
4. Formulas for standard errors of MLEs are well-known.  
Estimates of standard errors are available as by-product of numerical optimization (maximization) procedures.



## Case Study III: The Donner Party Example

## Case Study III: The Data

- ▶ Data:  $n=45$  pioneers

AGE	SEX	STATUS
23	MALE	DIED
40	FEMALE	SURVIVED
40	MALE	SURVIVED
30	MALE	DIED
28	MALE	DIED
40	MALE	DIED
...		

- ▶ AGE: Adults, 15-65 yrs old
- ▶ SEX: 15 Females, 30 Males
- ▶ BINARY OUTCOME: 25 Died, 20 Survived

- ▶ Questions: What are the odds of survival for a 20-yr old female? Compare the odds of survival to that of a male of the same age.

↓  
for females

## Case Study III: Binary Logistic Regression Additive Model

$$\log \left( \frac{\pi_i}{1 - \pi_i} \right) = \beta_0 + \beta_1 \text{Age}_{i1} + \beta_2 \text{Sex}_{i2}, \quad i = 1, \dots, 45$$

$(15, 65)$        $(0, 1)$   
↓                      ↓

$$0 < \pi < 1$$

- ✗ ▶ Cannot predict survival ( $\pi = 1$ ) or death ( $\pi = 0$ ) of a pioneer
- Descriptive { ▶ Can estimate:
  - ▶  $\pi_i$  (the probability of survival)
  - ▶ odds of survival and
  - ▶ log-odds of survival based on Age and Sex of a pioneer
- Inferential { ▶ Can be used to get point and interval estimates of odds ratios
- ▶ Can test which predictors are relevant to determine odds of survival

## Using R for fitting GLMs

- ▶ fitting function:

`glm(formula, family, data)`

- ▶ family: link function, distribution of  $Y$ .

Examples include binomial, gaussian, poisson, Gamma

- ▶ complementary functions:

- ▶ `coefficients()`: coefficient estimates
- ▶ `summary()`: prints a summary of results
- ▶ `anova()`: produces an analysis of variance table
- ▶ residuals
- ▶ deviance

- ▶ Optimization technique: Fisher Scoring / IWLS

### Case Study III: Fitted equations

Using defaults,  $\pi = P(SURVIVED)$ :

$$(1) \quad \log \left( \frac{\hat{\pi}_i}{1 - \hat{\pi}_i} \right) = 3.23 - 0.078 \text{Age}_i - 1.60 \mathbb{1}_{\text{Male},i}$$

Using other reference status,  $\pi = P(DIED)$ :

$$\log \left( \frac{\hat{\pi}_i}{1 - \hat{\pi}_i} \right) = -3.23 + 0.078 \text{Age}_i + 1.60 \mathbb{1}_{\text{Male},i}$$

$x_2 = 1$  if Males

$$\log \left( \frac{\pi_i}{1 - \pi_i} \right) = - \log \left( \frac{1 - \pi_i}{\pi_i} \right) = \log \left( \frac{\pi_i}{1 - \pi_i} \right)^{-1}$$

Using sex reference group as Males,  $\pi = P(SURVIVED)$ :

$$(3) \quad \log \left( \frac{\hat{\pi}_i}{1 - \hat{\pi}_i} \right) = 1.63 - 0.078 \text{Age}_i + 1.60 \mathbb{1}_{\text{Female},i}$$

Using sex reference group as Males,  $\pi = P(DIED)$ :

$$\log \left( \frac{\hat{\pi}_i}{1 - \hat{\pi}_i} \right) = -1.63 + 0.078 \text{Age}_i - 1.60 \mathbb{1}_{\text{Female},i}$$

$x_2 = 1$  if Female

## Case Study III: Using Fitted equation

Using the fitted equation for  $\pi = P(SURVIVED)$ :

$$\log \left( \frac{\hat{\pi}_i}{1 - \hat{\pi}_i} \right) = 1.63 - 0.078 \text{Age}_i + 1.60 \mathbb{1}_{\text{Female},i}$$

Q: Estimate the log odds, odds and probability of survival for a:

	Log odds, $\log\left(\frac{\hat{\pi}}{1-\hat{\pi}}\right)$	Odds, $\frac{\hat{\pi}}{1-\hat{\pi}}$	$\hat{\pi}$
(i) 20-yr old Female	$1.63 - 0.078(20) + 1.6(1)$		
(ii) 40-yr old Female	$1.63 - 0.078(40) + 1.6$		
(iii) 20-yr old Male	$1.63 - 0.078(20)$		
(iv) 40-yr old Male	$1.63 - 0.078(40)$		

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

$$\frac{\text{odds}}{1 + \text{odds}}$$

$$e^{\log \text{odds}}$$

$$1_F = 1$$

$$1_F = 0$$

## Case Study III: Using Fitted equation

Using the fitted equation for  $\pi = P(SURVIVED)$ :

$$\log \left( \frac{\hat{\pi}_i}{1 - \hat{\pi}_i} \right) = 1.63 - 0.078Age_i + 1.60\mathbb{1}_{Female,i},$$

Q: Estimate the log odds, odds and probability of survival for a:

	Log odds, $\log(\frac{\hat{\pi}}{1-\hat{\pi}})$	Odds, $\frac{\hat{\pi}}{1-\hat{\pi}}$	$\hat{\pi}$
(i) 20-yr old Female	1.67	5.31	0.84
(ii) 40-yr old Female	0.11	1.12	0.53
(iii) 20-yr old Male	0.07	1.07	0.52
(iv) 40-yr old Male	-1.49	0.225	0.18

Qs: Compare the odds of survival for a 40-yr old Female to that of a 20-yr old Female. Compare the odds of survival for a 20-yr old Female to that of a Male of the same age.



## Case Study III: Odds Ratios

1. Compare the odds of survival for a 40-yr old Female to that of a 20-yr old Female.

$$\frac{1.12}{5.31} = 0.21 \approx \frac{1}{5}$$

Hence, the odds of survival for a 20-yr old Female are about 5 times the odds for a 40-yr old Female.

2. Compare the odds of survival for a 20-yr old Female to that of a Male of the same age.

$$\frac{5.31}{1.07} = 4.96 \approx 5$$

$$\frac{1.12}{0.225} = 4.96$$

Hence, the odds of survival for a 20-yr old Female are about 5 times the odds for a Male of the same age.

## Interpreting coefficients of a Binary Logistic model

For  $\pi = P(Y = 1)$  <sup>"Success"</sup> we model

$$\log \left( \frac{\pi}{1 - \pi} \right) = \log \text{odds}_{\{Y=1\}} = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$$

Let  $\omega$  be the odds that  $Y=1$  based on  $X_1, \dots, X_p$ , then

$$\omega = \exp\{\beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p\}.$$

**Interpretation of  $\beta_1$ :** Holding  $X_2, \dots, X_p$  fixed, the ratio of the odds ('ODDS RATIO') that  $Y=1$  at  $X_1=a$  to  $X_1=b$  is

$$\frac{\omega_a}{\omega_b} = \exp\{\beta_1(a - b)\} = \frac{e^{\beta_0 + \beta_1 a + \beta_2 X_2 + \cdots + \beta_p X_p}}{e^{\beta_0 + \beta_1 b + \beta_2 X_2 + \cdots + \beta_p X_p}}$$

\* If  $X_1$  increases by 1 unit, holding all other  $X$ 's constant, the odds that  $Y=1$  change by a multiplicative factor of  $e^{\beta_1}$ .

$$e^{2\beta_1} = e^{\beta_1} * e^{\beta_1}$$

## Case Study III: Using coefficients to find Odds Ratios

$$\log\left(\frac{\hat{\pi}}{1-\hat{\pi}}\right) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 = 1.63 - 0.078\text{Age} + 1.60\mathbb{1}_{\text{Female}}$$

1. (Fixed Sex) Compare the odds of survival for a 40-yr old (Female/Male) to that of a 20-yr old (Female/Male).

$$\exp\{-0.78(40 - 20)\} = 0.21 \approx \frac{1}{5}$$

Hence, the odds of survival for a 20-yr old are about 5 times the odds for a 40-yr old of the same sex.

2. (Fixed Age) Compare the odds of survival for a 20-yr old Female to that of a Male of the same age.

$$\exp\{1.60(1 - 0)\} = 4.95 \approx 5$$

Hence, the odds of survival for a Female are about 5 times the odds for a Male of the same age.

## Next

- ▶ Confidence interval for Odds Ratio
- ▶ Testing  $\beta$ 's  $\rightarrow$  Higher-order Models