# **Applied Logistic Regression**

#### Week 3

- 1. Homework week 2: highlights
- 2. Confidence interval for  $\beta$  and  $\pi$
- 3. Interpretation of coefficients
- 4. Dichotomous independent variable
- 5. Computer output: STATA
- 6. Homework

Stanley Lemeshow, Professor of Biostatistics College of Public Health, The Ohio State University



It follows from maximum likelihood theory that the estimator of the variance of the estimated coefficients is the inverse of the observed information matrix. The observed information matrix is the matrix of second partial derivatives evaluated at the MLE,  $\hat{\beta}$ .

In the case of a model with p covariates the elements in the matrix are:

$$\hat{I}\left(\hat{\boldsymbol{\beta}}\right) = \begin{bmatrix} -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}} = \sum_{i=1}^{n} \tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) & -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}\partial \beta_{1}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) & -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) \\ -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}\partial \beta_{1}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) & -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{1}^{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}^{2}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) & \cdots & -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) \\ -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) & -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{1}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) \\ -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) & -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{1}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\boldsymbol{x}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) \\ -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) & -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{1}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\boldsymbol{x}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) \\ -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) & -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{1}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\boldsymbol{x}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) \\ -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\left(1-\tilde{\boldsymbol{x}}\left(\boldsymbol{x}_{i}\right)\right) \\ -\frac{\partial^{2} L\left(\boldsymbol{\beta}\right)}{\partial \beta_{0}^{2}\partial \beta_{2}} = \sum_{i=1}^{n} \boldsymbol{x}_{i}\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x$$

$$\begin{split} &\text{if we let} \qquad X = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{bmatrix}_{n \times (p+1)} \\ &\hat{V} = \begin{bmatrix} \hat{\pi}\left(x_1\right)\left(1-\hat{\pi}\left(x_1\right)\right) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \hat{\pi}\left(x_1\right)\left(1-\hat{\pi}\left(x_1\right)\right) \end{bmatrix}_{n \times n} \\ &= \operatorname{diag}\left(\hat{\pi}\left(x_i\right)\left(1-\hat{\pi}\left(x_i\right)\right)\right) \\ &\text{then} \qquad \hat{I}\left(\hat{\beta}\right) = \begin{bmatrix} \mathbf{X'VX} \end{bmatrix}_{(p+1) \times (p+1)} \end{aligned}$$

It follows that the estimate of the covariance matrix of the estimated coefficients is the inverse of the estimated information matrix  $(2) [x^2x^2]^{-1}$ 

 $\widehat{Var}(\hat{\beta}) = [X'\hat{V}X]^{-1}$ 

This matrix is easily obtained from most packages.

It follows that the confidence interval estimator of a coefficient is

$$\widehat{\boldsymbol{\beta}}_{j} \pm \boldsymbol{Z}_{1-\alpha/2} \widehat{\boldsymbol{SE}} \left( \widehat{\boldsymbol{\beta}}_{j} \right)$$

where 
$$\widehat{SE}(\hat{\beta}_j) = \sqrt{\widehat{Var}(\hat{\beta}_j)}$$

Confidence interval for the logit for a single subject

The estimated logit is 
$$\hat{g}(x) = \sum_{j=0}^{p} \hat{\beta}_{j} x_{j}$$

and the estimate of its variance is 
$$\widehat{Var}(\hat{g}(x)) = x' \widehat{Var}(\hat{\beta}) x$$

Note: If we choose x as one of the subjects then we can obtain this from a computer package, otherwise we can use lincom.

# Hence the confidence interval for the logit is

$$\hat{g}(x) \pm z_{1-\alpha/2} \widehat{SE}(\hat{g}(x))$$

We use this to obtain a confidence interval for the fitted value or estimated logistic probability as follows:

$$\frac{e^{\hat{g}(x)\pm z_{1-\alpha/2}\widehat{SE}(\hat{g}(x))}}{1+e^{\hat{g}(x)\pm z_{1-\alpha/2}\widehat{SE}(\hat{g}(x))}}$$

# As an example consider a woman who has: lwt = 100 and race = black.

Hence the estimated probability and the end points of its confidence interval are:

$$\hat{\pi}$$
 (lwt = 100, Race = black) =  $\frac{e^{.3645094}}{1 + e^{.3645094}} = 0.590$ 

and

$$\frac{e^{-.596095}}{1+e^{-.596095}} = 0.355 \quad \text{and} \quad \frac{e^{1.325114}}{1+e^{1.325114}} = 0.790$$

## Alternatively, we could use the following command:

Question: What is the interpretation of this point estimate and confidence interval?

- we focus on the estimated coefficients for the independent variables in the model.

The coefficients represent the slope or rate of change of a function of the dependent variable per unit of change in the independent variable.

# Interpretation is then concerned with

- Determination of the functional relationship between independent and dependent variables
- Definition of unit of change for the independent variable.

# What function of y yields a linear function of the independent variable?

# [this function is known as the "link function"]

 In linear regression it is the identity function (i.e., y = y) since, by definition, y is linear in the parameters.

$$\mathbf{y} = \beta_0 + \beta_1 \mathbf{x}$$

In logistic regression, the link function is

$$g(x) = \ln \left\{ \frac{\pi(x)}{1 - \pi(x)} \right\} = \beta_0 + \beta_1 x$$

# How do we define the slope?

• In linear regression  $\beta_1 = \frac{\Delta y}{\Delta x}$ 

i.e., it is the difference  $y\Big|_{x=x+1}-y\Big|_{x=x}$  for any value of x.

e.g., let 
$$y(x) = \beta_0 + \beta_1 x$$
  
then  $y(x+1) = \beta_0 + \beta_1 (x+1) = \beta_0 + \beta_1 x + \beta_1$   
so  
 $y(x+1) - y(x) = \beta_0 + \beta_1 x + \beta_1 - (\beta_0 + \beta_1 x) = \beta_1$ 

Hence  $\beta_1$  = change in y corresponding to a unit change in x.

# e.g., Study of adolescent aged males

Suppose we find 
$$\hat{y} = \hat{\beta}_0 + 5x$$

 $\Rightarrow$   $\Delta$  in height of 1 inch is associated with a  $\Delta$  in weight of 5 pounds

# for logistic regression

$$g(x) = \ln\left\{\frac{\pi(x)}{1 - \pi(x)}\right\} = \beta_0 + \beta_1 x$$

$$g(x+1) = \ln\left\{\frac{\pi(x+1)}{1 - \pi(x+1)}\right\} = \beta_0 + \beta_1(x+1) = \beta_0 + \beta_1 x + \beta_1$$

SO

$$g(x+1)-g(x)=(\beta_0+\beta_1x+\beta_1)-(\beta_0+\beta_1x)=\beta_1$$

i.e.,  $\beta_1 = \Delta$  in logit for a 1 unit  $\Delta$  in x.

To interpret the coefficient,  $\beta_1$ , in a logistic regression model, we must place meaning on the difference between 2 logits.

- this will depend on the nature of the independent variable

- i.e., is it dichotomous
  - polychotomous
  - continuous

- we now consider each of these in turn.

#### We assume x is coded as 0,1

$$\Pr\left(y=1\big|X\right)=\pi\left(X\right)=\frac{e^{\beta_0+\beta_1X}}{1+e^{\beta_0+\beta_1X}}$$

# Then, under the logistic model we have

# Dependent variable (x) $x = 1 \qquad x = 0$ $\pi(1) = \frac{e^{\beta_0 + \beta_1}}{1 + e^{\beta_0 + \beta_1}} \qquad \pi(0) = \frac{e^{\beta_0}}{1 + e^{\beta_0}}$ $y = 0 \qquad 1 - \pi(1) = \frac{1}{1 + e^{\beta_0 + \beta_1}} \qquad 1 - \pi(0) = \frac{1}{1 + e^{\beta_0}}$ Total $1.0 \qquad 1.0$

The odds of the outcome being present among individuals with x = 1 is defined as

$$\frac{\Pr(y=1|x=1)}{\Pr(y=0|x=1)} = \frac{\pi(1)}{1-\pi(1)}$$

Similarly, the odds of the outcome being present among Individuals with x = 0 is

$$\frac{\Pr(y=1|x=0)}{\Pr(y=0|x=0)} = \frac{\pi(0)}{1-\pi(0)}$$

#### **RECALL**:

We have defined the "logit" as

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

Hence, the logit ≡ log of the odds or "log odds"

Specifically,

$$g(1) = \ln \begin{bmatrix} \pi(1) \\ 1 - \pi(1) \end{bmatrix}$$

$$g(0) = \ln \begin{bmatrix} \pi(0) \\ 1 - \pi(0) \end{bmatrix}$$

# The "odds ratio" is defined as the ratio of the odds for x = 1 to the odds for x = 0

i.e., 
$$\frac{\pi(1)}{1-\pi(1)}$$
 
$$OR = \frac{1-\pi(1)}{\pi(0)}$$
 
$$1-\pi(0)$$

The log of the odds ratio, termed "log-odds ratio" is

$$\ln(OR) = \ln \left[ \frac{\pi(1)}{1 - \pi(1)} \right] = g(1) - g(0) \Leftarrow \text{logit difference}$$

$$1 - \pi(0)$$

Now:

$$OR = \frac{\left[\frac{e^{\beta_0 + \beta_1}}{1 + e^{\beta_0 + \beta_1}}\right] \left[\frac{1}{1 + e^{\beta_0}}\right]}{\left[\frac{e^{\beta_0}}{1 + e^{\beta_0}}\right] \left[\frac{1}{1 + e^{\beta_0 + \beta_1}}\right]} = \frac{e^{\beta_0 + \beta_1}}{e^{\beta_0}} = e^{\beta_1}$$

Hence, for a dichotomous independent variable

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

and the logit difference, or log odds ratio is

$$\ln(OR) = \ln(e^{\beta_1}) = \beta_1$$

\*\* This is why logistic regression has proven such a powerful tool in epidemiologic research.

## e.g., let us consider again the CHD/AGE data.

#### We create a new variable as

let 
$$x = \begin{cases} 1 \text{ if } AGE \ge 55 \text{ years} \\ 0 \text{ if } AGE < 55 \text{ years} \end{cases}$$

## This results in the following $2 \times 2$ table:

		AGE			
		≥ 55(1)	< 55(0)		
CHD	1	21	22	43	
	0	6	51	<b>57</b>	
		<b>27</b>	<b>73</b>	100	

#### The likelihood for these data is

$$I\left(\underline{\beta}\right) = \pi \left(1\right)^{21} \left[1 - \pi \left(1\right)\right]^{6} \pi \left(0\right)^{22} \left[1 - \pi \left(0\right)\right]^{51}$$

Use of a logistic regression program provides the following:

VARIABLE	β̂	$\widehat{SE}(\hat{\beta})$	$\widehat{\widehat{SE}}(\widehat{\widehat{\beta}})$	<b>O</b> R
AGE	2.0935	0.5285	3.961	8.1
CONSTANT	-0.8408	0.2551	- 3. 296	

The value of 
$$\widehat{OR} = e^{2.0935} = 8.1$$

NOTE: we could have computed OR directly as:

$$\widehat{OR} = \frac{21 \times 51}{22 \times 6} = 8.11$$

Also: 
$$\hat{\beta}_1 = \ln(\widehat{OR}) = \ln(8.11) = 2.0935$$

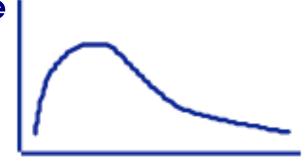
This will be an enormously useful fact when there is more than 1 independent variable where

$$e^{\hat{\beta}_i}$$
 = adjusted  $OR$ 

where the adjustment is for all other independent variables currently in the model.

# **Confidence Interval Estimation**

Because the distribution of OR tends to be skewed to the right, it is clearly not normally distributed. Thus, confidence intervals are usually based on  $\hat{\beta}_1$ 



(since In(OR) is more nearly normally distributed).

For sufficiently large samples

$$\hat{\beta}_{1} \sim N(\beta_{1}, Var(\hat{\beta}_{1}))$$

We can estimate the variance in the case when x is dichotomous as

$$\widehat{Var}\left(\hat{\beta}_{1}\right) = \left[\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}\right] \quad \text{and} \quad \widehat{SE}\left(\hat{\beta}_{1}\right) = \sqrt{\widehat{Var}\left(\hat{\beta}_{1}\right)}$$

cell frequencies in the  $2 \times 2$  table of  $y \times x$ 

# Thus, we compute

$$e^{\hat{eta}_{1}\pm z_{1-lpha/2}\widehat{SE}\left(\hat{eta}_{1}
ight)}$$

# **Example:**

$$\hat{\beta}_1 = 2.0935$$
  $\widehat{SE}(\hat{\beta}_1) = \left[\frac{1}{21} + \frac{1}{22} + \frac{1}{6} + \frac{1}{51}\right]^{\frac{1}{2}} = 0.5285$ 

and a 95% Cl is

$$2.88 \le OR \le 22.86$$

# An alternative method of coding design variables for a Dichotomous risk factor is as follows:

If we use the deviation from means coding with a dichotomous variable then we have X D

Can we obtain the odds ratio with this coding?

$$\ln(OR) = g(x = 1) - g(x = 0)$$

$$= g(D = +1) - g(D = -1)$$

$$= (\beta_0 + \beta_1) - (\beta_0 - \beta_1) = 2\beta_1$$

$$\therefore \widehat{OR} = e^{2\hat{\beta}_1}$$

When we use deviation from means coding the confidence

interval is:  $e^{2\hat{\beta}_1 \pm 1.96 \times 2 \times \widehat{SE}(\hat{\beta}_1)}$ 

#### . gen aged=age>=55

#### . logit chd aged

Iteration 0: log likelihood = -68.331491
Iteration 1: log likelihood = -59.020453
Iteration 2: log likelihood = -58.979594
Iteration 3: log likelihood = -58.979565

Logit estimates	Number of obs	=	100
	LR chi2(1)	=	18.70
	Prob > chi2	=	0.0000
$Log\ likelihood = -58.979565$	Pseudo R2	=	0.1369

chd		Std. Err.	P> z	[95% Conf.	Interval]
aged	2.093546 8407832	.5285335	0.000	1.057639 -1.340718	3.129453 3408487

#### . logistic chd aged

Logit estimates				Number	; =	100	
				LR chi2	2(1)	=	18.70
				Prob >	chi2	=	0.0000
$Log\ likelihood = -58.979565$				Pseudo	R2 =	=	0.1369
	l   Odds Ratio		z	P> z	[95%	Conf.	Interval]
	8.113636		3.961		2.879	566	22.86147