



JOHNS HOPKINS  
BLOOMBERG SCHOOL  
of PUBLIC HEALTH

## Lecture 3

### Assessing confounding in logistic regression models MLE in logistic regression models

# Review of Lecture 2

- ▶ Regression adjustment in logistic regression
- ▶ Model B:
- ▶ Model D, binary age:
- ▶ Model D, continuous age:



# Model D: Parameter interpretation and estimation

```
modelD = glm(bigexp~mscd+older,data=data1,family="binomial")  
summary(modelD)$coeff
```

##	Estimate	Std. Error	z value	Pr(> z )
## (Intercept)	-0.9577826	0.02700779	-35.46321	1.815505e-275
## mscd	1.6549130	0.06803662	24.32386	1.096494e-130
## older	0.5638298	0.04104938	13.73540	6.230701e-43

```
lincom(modelD,c("mscd","older"),eform=TRUE)
```

##	Estimate	2.5 %	97.5 %	Chisq	Pr(>Chisq)
## mscd	5.232625	4.57938	5.979054	591.65	1.096494e-130
## older	1.75739	1.621537	1.904625	188.6613	6.230701e-43

You practice: Use the output above, interpret  $\exp(\hat{\beta}_2)$ .

# Model D: Adjustment for continuous covariates

```
modelDagecont = glm(bigexp~mscd+lastage,data=data1,family="binomial")
summary(modelDagecont)$coeff
```

	Estimate	Std. Error	z value	Pr(> z )
## (Intercept)	-2.27990966	0.099135981	-22.99780	4.903428e-117
## mscd	1.60502065	0.068269770	23.50998	3.224831e-122
## lastage	0.02574057	0.001599682	16.09105	2.947835e-58

```
lincom(modelDagecont,c("mscd","lastage"),eform=TRUE)
```

	Estimate	2.5 %	97.5 %	Chisq	Pr(>Chisq)
## mscd	4.977962	4.35452	5.690664	552.719	3.224831e-122
## lastage	1.026075	1.022863	1.029297	258.922	2.947835e-58

- Interpret both of the coefficients:



# Assessing confounding in logistic regression

- ▶ Question: Is age a confounder for the big expenditure vs. MSCD relationship?
  - ▶ Is age associated with MSCD?
  - ▶ Is age associated with having a big expenditure?
  - ▶ Is age in the causal pathway between MSCD and having a big expenditure?
- ▶ We can answer questions 1 and 2 using statistical analyses
  - ▶ Question 3 is not a statistical question

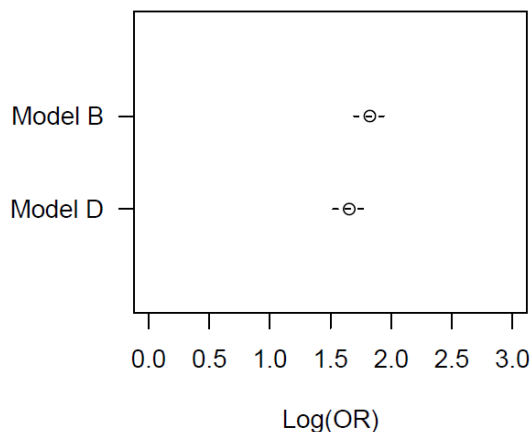


# Assessing confounding in logistic regression

From the analysis of the NMES data, we found:

```
##      Estimate    2.5 %    97.5 %    Chisq    Pr(>Chisq)
## mscd 1.825045  1.694177  1.955913  747.095  1.718138e-164
```

```
##      Estimate    2.5 %    97.5 %    Chisq    Pr(>Chisq)
## mscd 1.654913  1.521564  1.788262  591.65   1.096494e-130
```



Estimated difference:  $1.82 - 1.65 = 0.175$   
95% CI for the difference: 0.14 to 0.20

What do you think?

# Assessing confounding in logistic regression

- ▶ You can use the process described above for general linear model or generalized linear models with log links
  - ▶ However, it gets tricky for other link functions, e.g. the logit link.
- ▶ For rest of the lecture consider:
  - ▶  $Y$  = binary outcome
  - ▶  $X$  = primary binary exposure variable
  - ▶  $Z$  = potential confounding variable

$$\text{Model M: } \textit{logit}[Pr(Y = 1|X)] = \beta_{0m} + \beta_{1m}X$$

$$\text{Model C: } \textit{logit}[Pr(Y = 1|X, Z)] = \beta_{0c} + \beta_{1c}X + \beta_{2c}Z$$



# Non-linearity effect in logistic regression models

- ▶ Assume X and Z are independent, i.e. no confounding
- ▶ You can show that  $|\beta_{1c}| > |\beta_{1m}|$  and the difference depends on the relationship between X and Z on Y and the variance of Z.
- ▶ Difference is referred to as the “non-linearity effect”
- ▶ This feature of the logistic regression model is known as “non-collapsibility”





# Non-linearity effect in logistic regression models

- ▶ Implications for evaluating confounding:
- ▶  $|\beta_{1c}| < |\beta_{1m}|$ 
  - ▶ You have identified “positive confounding” despite the non-linearity effect
- ▶  $|\beta_{1c}| > |\beta_{1m}|$ 
  - ▶ You may have “negative confounding” or you may be observing the non-linearity effect
- ▶  $\beta_{1c}$  and  $\beta_{1m}$  have different signs!
  - ▶ You may have “qualitative confounding” or you may be observing the non-linearity effect



# Linear models: no non-linearity effect

- ▶  $Y \sim \text{Normal}$ ,  $X$  and  $Z$  independent

- ▶ Marginal model:

- ▶ Conditional Model:

- ▶ Marginal model coefficient:

$$\begin{aligned} E(Y|X=1) - E(Y|X=0) &= E_Z[E(Y|X=1, Z) - E(Y|X=0, Z)] \\ &= E_Z[(\beta_{0c} + \beta_{1c} + \beta_{2c}Z) - (\beta_{0c} + \beta_{2c}Z)] \\ &= \beta_{1c} \end{aligned}$$



# Logistic regression: non-linearity effect

- ▶ Consider the marginal and conditional odds ratios

$$\exp(\beta_{1m}) = \frac{\exp(\beta_{0m} + \beta_{1m})}{\exp(\beta_{0m})} = \frac{Pr(Y = 1|X = 1)/Pr(Y = 0|X = 1)}{Pr(Y = 1|X = 0)/Pr(Y = 0|X = 0)}$$

$$\exp(\beta_{1c}) = \frac{\exp(\beta_{0c} + \beta_{1c} + \beta_{2c}Z)}{\exp(\beta_{0m} + \beta_{2c}Z)} = \frac{Pr(Y = 1|X = 1, Z)/Pr(Y = 0|X = 1, Z)}{Pr(Y = 1|X = 0, Z)/Pr(Y = 0|X = 0, Z)}$$

- ▶ When would these be the same?
  - ▶  $\beta_{2c} = 0$ , Y and Z independent
  - ▶  $\beta_{1c} = 0, \beta_{1m} = 0$ , Y and X independent
  - ▶  $\text{Var}(Z) = 0$

- ▶ Why aren't they the same?

$$E(Y|X) = Pr(Y = 1|X) = E_Z\left[\frac{\exp(\beta_{0c} + \beta_{1c}X + \beta_{2c}Z)}{1 + \exp(\beta_{0c} + \beta_{1c}X + \beta_{2c}Z)}\right]$$

# Simulation: Non-linearity effect

- ▶ Assume the following model:

$$\text{Logit}[\Pr(Y = 1|X, Z)] = -2 + 0.4 X + Z$$

where  $Z \sim N(0, 2)$ , and  $X$  and  $Z$  are independent

- ▶ This model says that regardless of the value of  $Z$ , the relative odds of  $Y = 1$  comparing persons with  $X = 1$  to persons with  $X = 0$  are  $\exp(0.4) = 1.5$



## Simulation: Non-linearity effect

This model says that regardless of the value of Z, the relative odds of  $Y = 1$  comparing persons with  $X = 1$  to persons with  $X = 0$  are  $\exp(0.4) = 1.5$

Consider a person with  $Z = 0$ :

$$Pr(Y = 1|X = 1, Z = 0) = \frac{\exp(-2 + 0.4)}{1 + \exp(-2 + 0.4)} = 0.17$$

$$Pr(Y = 1|X = 0, Z = 0) = \frac{\exp(-2)}{1 + \exp(-2)} = 0.12$$

$$OR(Y, X|Z = 0) = \frac{0.17/0.83}{0.12/0.88} = 1.5$$



## Simulation: Non-linearity effect

This model says that regardless of the value of  $Z$ , the relative odds of  $Y = 1$  comparing persons with  $X = 1$  to persons with  $X = 0$  are  $\exp(0.4) = 1.5$

Consider a person with  $Z = 2$ :

$$Pr(Y = 1|X = 1, Z = 2) = \frac{\exp(-2 + 0.4 + 2)}{1 + \exp(-2 + 0.4 + 2)} = 0.60$$

$$Pr(Y = 1|X = 0, Z = 2) = \frac{\exp(-2 + 2)}{1 + \exp(-2 + 2)} = 0.5$$

$$OR(Y, X|Z = 2) = \frac{0.6/0.4}{0.5/0.5} = 1.5$$



# Simulation: Non-linearity effect

What about the marginal probabilities?

$$Pr(Y = 1|X) = \int_z Pr(Y = 1|X, z)f(z)dz$$

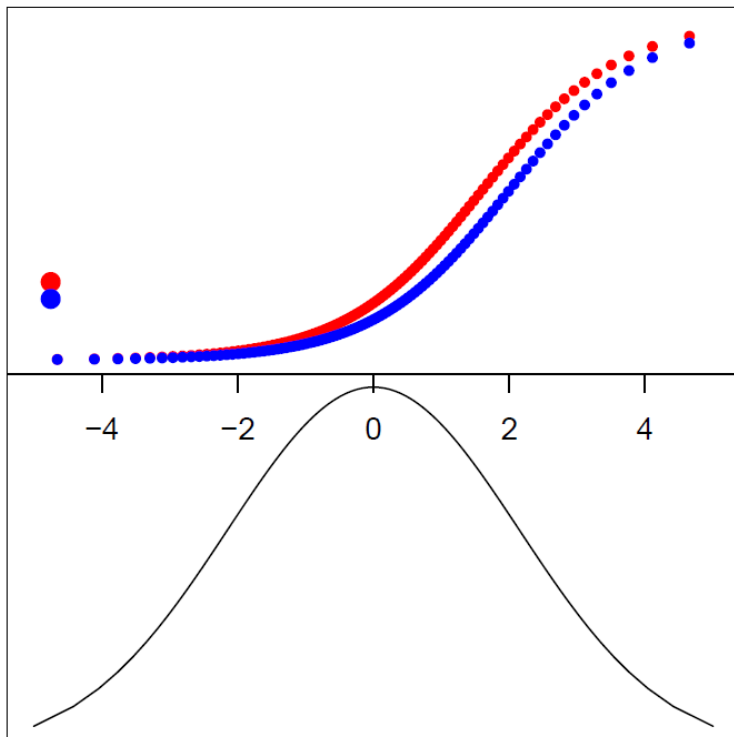
The marginal probabilities are a weighted average of the conditional probabilities with weights determined by the normal density

$$Pr(Y = 1|X = 1) = 0.23$$

$$Pr(Y = 1|X = 0) = 0.18$$

$$OR(Y, X) = \frac{0.23/0.77}{0.18/0.82} = 1.36$$

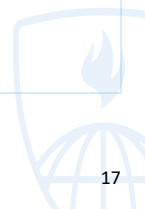
## Simulation: Non-linearity effect





# Very important note!

- ▶ The non-linearity effect is roughly the same for estimation of  $\beta_{1c}$  and  $\text{se}(\hat{\beta}_{1c})$ .
- ▶ So we would expect the Z statistics for  $\beta_{1c}$  and  $\beta_{1m}$  to be roughly the same if little to no confounding is present (e.g. X and Z are independent).
- ▶ So we would expect the Z statistics for  $\beta_{1c}$  and  $\beta_{1m}$  to be different if Z is a confounder.

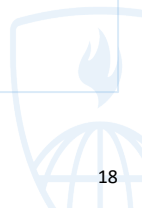


# Simulation: Confounding present

- ▶ Mimic the simulation described in Janes et al (Biostatistics, 2010).

Assume the following:

- We simulate 1000 samples of 250 persons, half with exposure ( $X = 1$ ) and half without ( $X = 0$ ).
  - We generate  $Z$  as follows:  $Z|X \sim N(\alpha_0 + \alpha_1 X, 1)$
  - We generate  $Y$  from:  $\text{logit}[Pr(Y = 1|X, Z)] = \beta_0 + \beta_1 X + \beta_2 Z$ .
  - We set  $\alpha_0 = 0$ ,  $\beta_0 = 0$  and  $\beta_1 = \log(2)$ .
- ▶ Simulation scenarios:
    - ▶ No confounding / no non-linearity:  $\alpha_1 \approx 0, \beta_1 \approx 0$
    - ▶ No confounding / non-linearity:  $\alpha_1 \approx 0, \beta_1$  large
    - ▶ “Small” confounding:  $\alpha_1 > 0, \beta_1 \approx 0$
    - ▶ Confounding:  $\alpha_1, \beta_1$  large



# Simulation: Confounding present

##		a1	b1	beta1m	beta1	beta1m-beta1	Z1m	Z1	Z1m-Z1
## 1		0.01	0.05	0.701	0.704	-0.003	2.653	2.653	0.000
## 2		0.01	1.50	0.494	0.689	-0.195	1.902	2.183	-0.281
## 3		1.00	0.05	0.752	0.711	0.041	2.839	2.389	0.450
## 4		1.00	1.50	1.617	0.707	0.910	5.312	1.951	3.361

- ▶ No confounding / no non-linearity:  $\alpha_1 \approx 0, \beta_1 \approx 0$ 
  - ▶ Coefficients and test statistics are the same
- ▶ No confounding / non-linearity:  $\alpha_1 \approx 0, \beta_1$  large
  - ▶ Conditional coefficients are different, test statistics are roughly the same
- ▶ “Small” confounding:  $\alpha_1 > 0, \beta_1 \approx 0$ 
  - ▶ Non-linearity effect is small
  - ▶ Marginal coefficient > conditional coefficient -> confounding
  - ▶ Test statistics differ
- ▶ Confounding:  $\alpha_1, \beta_1$  large
  - ▶ Marginal coefficient > conditional coefficient -> confounding
  - ▶ Test statistics differ

# Shifting gears to estimation! MLE in linear models

- Define the  $(p+1) \times 1$  vector of covariates for subject  $i$  as  $x_i = (1, x_{1i}, x_{2i}, \dots, x_{pi})$ .
- Define the  $(p+1) \times 1$  vector of association parameters  $\beta = (\beta_0, \beta_1, \dots, \beta_p)$ .

$$Y_i = \mu_i + \epsilon_i, \epsilon_i \text{ iid } N(0, \sigma^2)$$

$$E(Y_i) = \mu_i(\beta) = x_i' \beta$$

The score equation,  $U(\beta) = \frac{\partial \log L(\beta | y_i)}{\partial \beta} = \sum_{i=1}^n x_i (y_i - \mu_i(\beta))$ .

Setting  $U(\beta) = \sum_{i=1}^n x_i (y_i - \mu_i(\beta)) = 0$  and solving for  $\beta$  produced:

$$\hat{\beta} = (X'X)^{-1}X'Y$$

where  $X$  is the  $n \times p$  matrix of stacked row vectors  $x_i'$  and  $Y$  is the  $1 \times n$  vector of responses.

# MLE in logistic models

Assume the following model:

- $Y_i \sim \text{Bernoulli}(\mu_i)$  for  $i = 1, \dots, n$  independent observations.
- Define the vector of covariates for subject  $i$  as  $x_i = (1, x_{1i}, x_{2i}, \dots, x_{pi})$ .
- Define the vector of association parameters  $\beta = (\beta_0, \beta_1, \dots, \beta_p)$ .
- Assume the logit link such that:

$$\log\left(\frac{\mu_i}{1 - \mu_i}\right) = x_i^T \beta \rightarrow \mu_i = \frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}}$$

NOTE: We should really write  $\mu_i(x_i, \beta)$  i.e.  $\mu_i$  is a function of  $x_i$  and  $\beta$ . In this handout, I will simplify this to  $\mu_i(\beta)$ .



# MLE in logistic models

We can express the likelihood function as:

$$\begin{aligned}L(\beta|y) &= Pr(Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n|\beta) \\&= \prod_{i=1}^n Pr(Y_i = y_i|\beta) \\&= \prod_{i=1}^n \mu_i(\beta)^{y_i} [1 - \mu_i(\beta)]^{1-y_i}\end{aligned}$$

The log-likelihood function is:

$$\log[L(\beta|y)] = \sum_{i=1}^n y_i \log[\mu_i(\beta)] + (1 - y_i) \log[1 - \mu_i(\beta)]$$



# MLE in logistic models

The score equation,  $U(\beta)$  is the derivative of the log-likelihood function with respect to  $\beta$ .

$$\begin{aligned}U(\beta) &= \frac{\partial \log[L(\beta|y)]}{\partial \beta} \\&= \sum_{i=1}^n y_i \frac{\partial \log[\mu_i(\beta)]}{\partial \beta} + (1 - y_i) \frac{\partial \log[1 - \mu_i(\beta)]}{\partial \beta}\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial \beta} \log[\mu_i(\beta)] &= \frac{\partial}{\partial \beta} \log \left( \frac{e^{x_i' \beta}}{1 + e^{x_i' \beta}} \right) \\&= \frac{\partial}{\partial \beta} [x_i' \beta - \log(1 + e^{x_i' \beta})] \\&= x_i - x_i \frac{e^{x_i' \beta}}{1 + e^{x_i' \beta}} \\&= x_i [1 - \mu_i(\beta)]\end{aligned}$$

# MLE in logistic models

The score equation,  $U(\beta)$  is the derivative of the log-likelihood function with respect to  $\beta$ .

$$\begin{aligned}U(\beta) &= \frac{\partial \log[L(\beta|y)]}{\partial \beta} \\&= \sum_{i=1}^n y_i \frac{\partial \log[\mu_i(\beta)]}{\partial \beta} + (1 - y_i) \frac{\partial \log[1 - \mu_i(\beta)]}{\partial \beta}\end{aligned}$$

For the next derivation, note that:

$$\frac{\partial \log[\mu_i(\beta)]}{\partial \beta} = \frac{1}{\mu_i(\beta)} \frac{\partial \mu_i(\beta)}{\partial \beta} \rightarrow \frac{\partial \mu_i(\beta)}{\partial \beta} = \mu_i \frac{\partial \log[\mu_i(\beta)]}{\partial \beta}$$

$$\begin{aligned}\frac{\partial}{\partial \beta} \log[1 - \mu_i(\beta)] &= -\frac{1}{1 - \mu_i(\beta)} \frac{\partial \mu_i(\beta)}{\partial \beta} \\&= \frac{-\mu_i(\beta)}{1 - \mu_i(\beta)} x_i [1 - \mu_i(\beta)] \\&= -\mu_i(\beta) x_i\end{aligned}$$





# MLE in logistic models

The score equation,  $U(\beta)$  is the derivative of the log-likelihood function with respect to  $\beta$ .

$$\begin{aligned}U(\beta) &= \frac{\partial \log[L(\beta|y)]}{\partial \beta} \\&= \sum_{i=1}^n y_i \frac{\partial \log[\mu_i(\beta)]}{\partial \beta} + (1 - y_i) \frac{\partial \log[1 - \mu_i(\beta)]}{\partial \beta} \\&= \sum_{i=1}^n y_i (x_i [1 - \mu_i(\beta)]) + (1 - y_i) [-\mu_i(\beta) x_i] \\&= \sum_{i=1}^n x_i (y_i - y_i \mu_i(\beta) + (-\mu_i(\beta)) + y_i \mu_i(\beta)) \\&= \sum_{i=1}^n x_i (y_i - \mu_i(\beta)) \\&= X^t(Y - \mu(\beta))\end{aligned}$$

# MLE in logistic models

NOTE: We will also need to know  $U'(\beta) = \frac{\partial U(\beta)}{\partial \beta}$

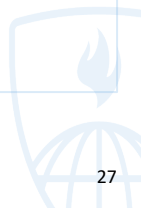
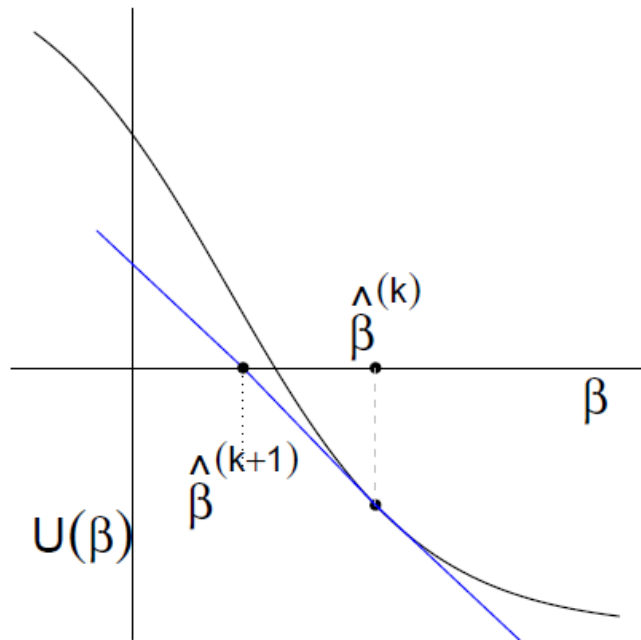
$$\begin{aligned}U'(\beta) &= \frac{\partial U(\beta)}{\partial \beta} \\&= \frac{\partial}{\partial \beta} X' (Y - \mu(\beta)) \\&= -X' \frac{\partial \mu_i(\beta)}{\partial \beta} \\&= -X' V X\end{aligned}$$

where we already showed that:

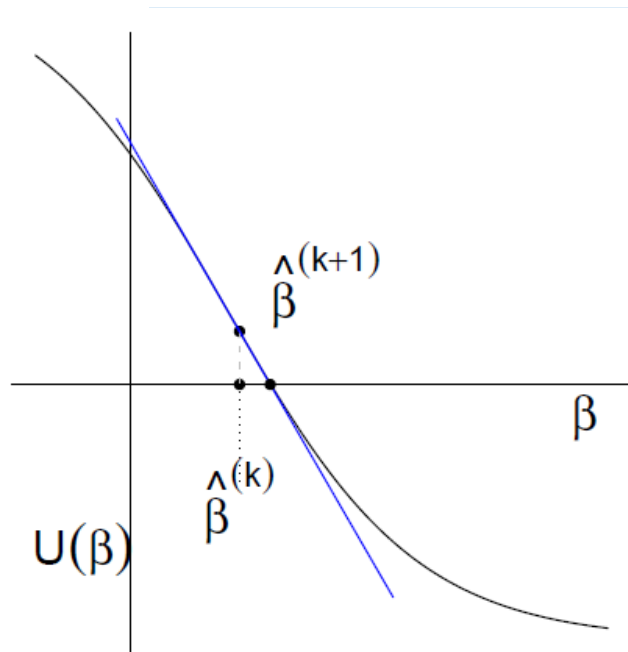
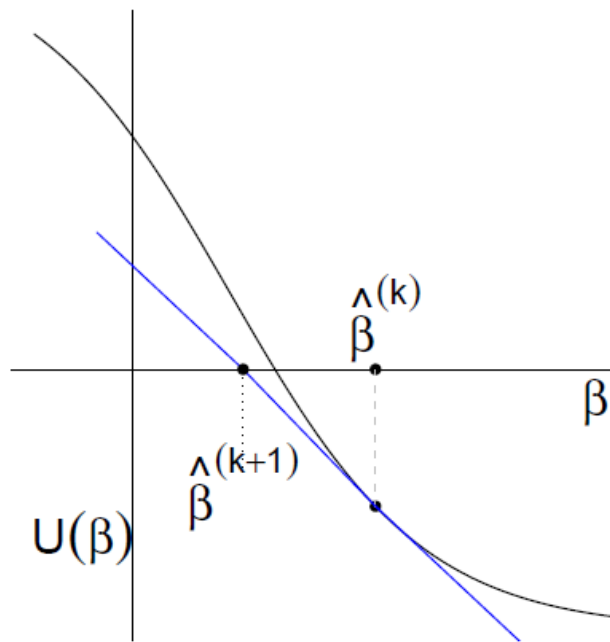
$$\frac{\partial \mu_i(\beta)}{\partial \beta} = \mu_i(\beta) \frac{\partial \log[\mu_i(\beta)]}{\partial \beta} = \mu_i(\beta)(1 - \mu_i(\beta))x_i$$

and  $V_{n \times n} = \text{diag}(\mu_i(\beta)[1 - \mu_i(\beta)])$ .

## Newton-Raphson Method to find “beta”



## Newton-Raphson Method to find “beta”



# Newton-Raphson Method to find “beta”

- Step 0: Pick an initial starting value for  $\beta$ , call this  $\hat{\beta}^{(k)}$ .
- Step 1: Compute the slope of  $U(\beta)$  at  $\hat{\beta}^{(k)}$ , i.e. compute  $U'(\hat{\beta}^{(k)})$ .
- Step 2: Construct the tangent line, which is a line that passes through the points  $(\hat{\beta}^{(k)}, U(\hat{\beta}^{(k)}))$  and  $(\hat{\beta}^{(k+1)}, 0)$  and has slope  $U'(\hat{\beta}^{(k)})$ .
- Step 3: Solve the following for  $\hat{\beta}^{(k+1)}$ :

$$U'(\hat{\beta}^{(k)}) = \frac{U(\hat{\beta}^{(k)}) - 0}{\hat{\beta}^{(k)} - \hat{\beta}^{(k+1)}}$$

$$[\hat{\beta}^{(k)} - \hat{\beta}^{(k+1)}]U'(\hat{\beta}^{(k)}) = U(\hat{\beta}^{(k)})$$

$$\hat{\beta}^{(k)} - \hat{\beta}^{(k+1)} = U'(\hat{\beta}^{(k)})^{-1}U(\hat{\beta}^{(k)})$$

$$\hat{\beta}^{(k+1)} = \hat{\beta}^{(k)} - U'(\hat{\beta}^{(k)})^{-1}U(\hat{\beta}^{(k)})$$

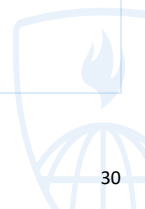
$$= U'(\hat{\beta}^{(k)})^{-1} \left( U'(\hat{\beta}^{(k)})\hat{\beta}^{(k)} - U(\hat{\beta}^{(k)}) \right)$$

- Step 4: Stop if  $|\hat{\beta}^{(k+1)} - \hat{\beta}^{(k)}|$  is small. If not, let  $k = k + 1$  and repeat Steps 2 through 4.

# Iteratively Re-weighted Least Squares

The general procedure is:

- Step 0: Set an initial value for  $\hat{\beta}^{(k)}$ ,  $k = 0$ .
- Step 1: Calculate:  $V^{(k)}$ ,  $\hat{\mu}(\hat{\beta}^{(k)})$ ,  $Z^{(k)}$ .
- Step 2: Update  $\hat{\beta}^{(k+1)} = (X^T V^{(k)} X)^{-1} (X^T V^{(k)} Z^{(k)})$
- Step 3: Stop if  $\sum_{j=1}^{p+1} \left( \hat{\beta}_j^{(k+1)} - \hat{\beta}_j^{(k)} \right)^2 < \epsilon$ ; if not, let  $k = k + 1$  and repeat Steps 2 and 3.



# Where to next?

- ▶ Inference within logistic regression models
  - ▶ Review -> some worked examples

