#### Lecture 19: Conditional Logistic Regression

Dipankar Bandyopadhyay, Ph.D.

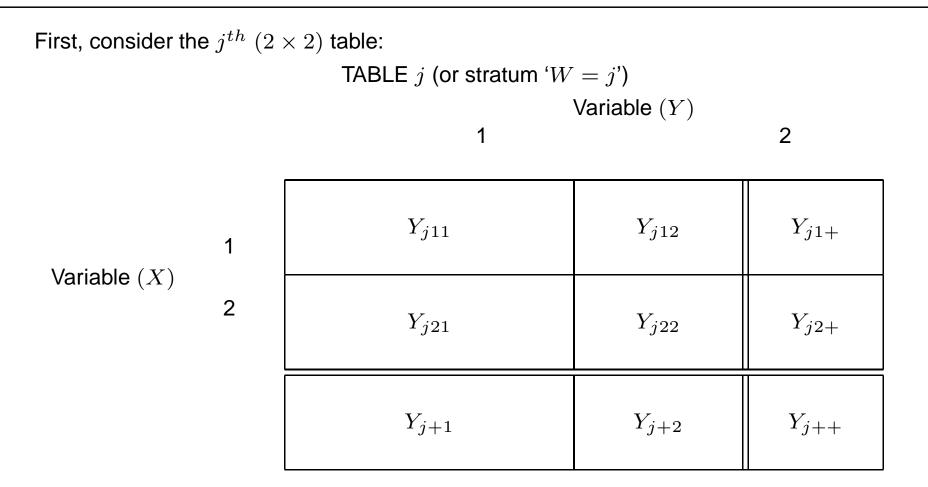
BMTRY 711: Analysis of Categorical Data Spring 2011 Division of Biostatistics and Epidemiology Medical University of South Carolina

#### Purpose

- 1. Eliminate unwanted nuisance parameters
- 2. Use with sparse data

Prior to the development of the conditional likelihood, lets review the unconditional (regular) likelihood associated with the logistic regression model.

- Suppose, we can group our covariates into *J* unique combinations
- and as such, we can form  $J(2 \times 2)$  tables



We can write the  $(2 \times 2)$  data table of probabilities for stratum W = j as

			Y	
		1	2	
X	1	$p_{j1}$	$1 - p_{j1}$	1
21	2	$p_{j2}$	$1 - p_{j2}$	1

### Unconditional Likelihood

- Let  $p_j$  be the probability that a subject in stratum j has a success.
- Then the probability of observing  $y_j$  events in stratum j is

$$P(Y_j = y_j) = \begin{pmatrix} n_j \\ y_j \end{pmatrix} p_j^{y_j} (1 - p_j)^{n_j - y_j}$$

- However, we know  $p_j$  is a function of covariates
- Without loss of generality, assume we are interested in two covariates,  $x_{j1}$  and  $x_{j2}$ , such that

$$p_j = \frac{e^{\beta_0 + \beta_1 x_{j1} + \beta_2 x_{j2}}}{1 + e^{\beta_0 + \beta_1 x_{j1} + \beta_2 x_{j2}}}$$

Then, the likelihood function can be written as

$$\begin{split} L(\vec{y}) &= \prod_{j=1}^{J} \binom{n_{j}}{y_{j}} p_{j}^{y_{j}} (1-p_{j})^{n_{j}-y_{j}} \\ &= \prod_{j=1}^{J} \left( \frac{e^{\beta_{0}+\beta_{1}x_{j1}+\beta_{2}x_{j2}}}{1+e^{\beta_{0}+\beta_{1}x_{j1}+\beta_{2}x_{j2}}} \right)^{y_{j}} \left( \frac{1}{1+e^{\beta_{0}+\beta_{1}x_{j1}+\beta_{2}x_{j2}}} \right)^{(n_{j}-y_{j})} \times \binom{n_{j}}{y_{j}} \\ &= \prod_{j=1}^{J} \frac{\left(e^{\beta_{0}+\beta_{1}x_{j1}+\beta_{2}x_{j2}}\right)^{y_{j}}}{(1+e^{\beta_{0}+\beta_{1}x_{j1}+\beta_{2}x_{j2}})^{n_{j}}} \times \binom{n_{j}}{y_{j}} \\ &= \frac{e^{\beta_{0}t_{0}+\beta_{1}t_{1}+\beta_{2}t_{2}}}{\prod_{j=1}^{J} \left(1+e^{\beta_{0}+\beta_{1}x_{j1}+\beta_{2}x_{j2}}\right)^{n_{j}}} \times \prod_{j=1}^{J} \binom{n_{j}}{y_{j}} \end{split}$$

where

$$t_k = \sum_{j=1}^J y_j x_{jk}$$

and

$$t_0 = \sum_{j=1}^J y_j$$

- The functions  $t_0$ ,  $t_1$ , and  $t_2$  are sufficient statistics for the data.
- Suppose we want to test  $\beta_2 = 0$  using a likelihood ratio test.
- Then

$$L(\beta_0, \beta_1) = \frac{e^{\beta_0 t_0 + \beta_1 t_1}}{\prod_{j=1}^{J} \left(1 + e^{\beta_0 + \beta_1 x_{j1}}\right)^{n_j}} \times \prod_{j=1}^{J} \left(\begin{array}{c} n_j \\ y_j \end{array}\right)$$

• and the LRT would equal

$$\begin{split} \Lambda(\vec{y}) &= -2(lnL(\beta_{0},\beta_{1}) - ln(L(\beta_{0},\beta_{1},\beta_{2}))) \\ &= -2ln\frac{L(\beta_{0},\beta_{1})}{L(\beta_{0},\beta_{1},\beta_{2})} \\ &= -2ln\frac{\frac{e^{\beta_{0}t_{0} + \beta_{1}t_{1}}}{\prod\limits_{j=1}^{J}\left(1 + e^{\beta_{0} + \beta_{1}x_{j1}}\right)^{n_{j}}} \times \prod\limits_{j=1}^{J}\left(\begin{array}{c}n_{j} \\ y_{j}\end{array}\right) \\ &= -2ln\frac{\frac{e^{\beta_{0}t_{0} + \beta_{1}t_{1} + \beta_{2}t_{2}}}{\prod\limits_{j=1}^{J}\left(1 + e^{\beta_{0} + \beta_{1}x_{j1} + \beta_{2}x_{j2}}\right)^{n_{j}}} \times \prod\limits_{j=1}^{J}\left(\begin{array}{c}n_{j} \\ y_{j}\end{array}\right) \end{split}$$

• Which is distributed **approximately** chi-square with 1 df.

- When you have small sample sizes, the chi-square approximation is not valid.
- We need a method to calculate the **exact** p-value of  $H_0: \beta_2 = 0$  from the exact null distribution of  $\Lambda(\vec{y})$ .
- Note that the distribution of  $\Lambda(\vec{y})$  depends on the exact distribution of  $\vec{y}$ .
- Thus, we need to derive the exact distribution of  $\vec{y}$ .

Under,  $H_0$ ,  $\beta_2 = 0$ , so we will begin with

$$P(\vec{y}) = L(\beta_0, \beta_1) = \frac{e^{\beta_0 t_0 + \beta_1 t_1}}{\prod_{j=1}^{J} \left(1 + e^{\beta_0 + \beta_1 x_{j1}}\right)^{n_j}} \times \prod_{j=1}^{J} \left(\begin{array}{c} n_j \\ y_j \end{array}\right)$$

- We cannot use the above probability statement because it relies on the population parameters (i.e., unknown constants)  $\beta_0$  and  $\beta_1$ .
- That is,  $\beta_0$  and  $\beta_1$  are a nuisance to our calculations and as such are called **nuisance parameters**
- Similar to what we did in the Fisher's Exact Test, we are going to condition out  $\beta_0$  and  $\beta_1$

Using Bayes' Law,

$$P(\vec{y}|t_0, t_1) = \frac{\frac{e^{\beta_0 t_0 + \beta_1 t_1}}{\prod\limits_{j=1}^{J} \left(1 + e^{\beta_0 + \beta_1 x_{j1}}\right)^{n_j}} \times \prod\limits_{j=1}^{J} \begin{pmatrix} n_j \\ y_j \end{pmatrix}}{\sum\limits_{\vec{y} \in \Gamma} \frac{e^{\beta_0 t_0 + \beta_1 t_1}}{\prod\limits_{j=1}^{J} \left(1 + e^{\beta_0 + \beta_1 x_{j1}}\right)^{n_j}} \times \prod\limits_{j=1}^{J} \begin{pmatrix} n_j \\ y_j \end{pmatrix}}$$

where  $\Gamma$  is the set of all vectors of y such that

$$\sum_{j=1}^{J} y_j = t_0$$
$$\sum_{j=1}^{J} y_j x_{1j} = t_1$$

and  $0 \le y_j \le n_j$  and  $y_j$  is an integer.

Then,  $P(\vec{y}|t_0, t_1)$  can be simplified to

$$P(\vec{y}|t_{0},t_{1}) = \frac{\frac{e^{\beta_{0}t_{0}+\beta_{1}t_{1}}}{\prod\limits_{j=1}^{J} \left(1+e^{\beta_{0}+\beta_{1}x_{j1}}\right)^{n_{j}}} \times \prod\limits_{j=1}^{J} \left(\begin{array}{c}n_{j}\\y_{j}\end{array}\right)}{\sum\limits_{\vec{y}\in\Gamma} \frac{e^{\beta_{0}t_{0}+\beta_{1}t_{1}}}{\prod\limits_{j=1}^{J} \left(1+e^{\beta_{0}+\beta_{1}x_{j1}}\right)^{n_{j}}} \times \prod\limits_{j=1}^{J} \left(\begin{array}{c}n_{j}\\y_{j}\end{array}\right)}{\sum\limits_{j=1}^{e^{\beta_{0}t_{0}+\beta_{1}t_{1}}} \prod\limits_{j=1}^{n_{j}} \left(1+e^{\beta_{0}+\beta_{1}x_{j1}}\right)^{n_{j}}} \times \prod\limits_{j=1}^{J} \left(\begin{array}{c}n_{j}\\y_{j}\end{array}\right)}{\sum\limits_{j=1}^{J} \left(1+e^{\beta_{0}+\beta_{1}x_{j1}}\right)^{n_{j}}} \times \sum\limits_{\vec{y}\in\Gamma} \prod\limits_{j=1}^{J} \left(\begin{array}{c}n_{j}\\y_{j}\end{array}\right)}{\sum\limits_{\vec{y}\in\Gamma} \prod\limits_{j=1}^{J} \left(\begin{array}{c}n_{j}\\y_{j}\end{array}\right)}{\sum\limits_{\vec{y}\in\Gamma} \prod\limits_{j=1}^{J} \left(\begin{array}{c}n_{j}\\y_{j}\end{array}\right)}}$$

which does not contain any unknown parameters and we can calculate the exact p-value.

Let  $\lambda$  by any value such that  $\Lambda(\vec{y}) = \lambda$  and denote

$$\Gamma_{\lambda} = \{ \vec{y} : \vec{y} \in \Gamma \text{ and } \Lambda(\vec{y}) = \lambda \}$$

Then

$$P(\Lambda(\vec{y}) = \lambda) = \sum_{\vec{y} \in \Gamma_{\lambda}} P(\vec{y}|t_0, t_1)$$

To obtain the exact p-value for testing  $\beta_2 = 0$  we need to:

- 1. Enumerate all values of  $\lambda$
- 2. Calculate  $P(\Lambda(\vec{y}) = \lambda)$

3. *p*-value equals the sum of the as extreme or more extreme values of  $P(\Lambda(\vec{y}) = \lambda)$ 

Or use SAS PROC LOGISTIC or LogXact

 Effectiveness of immediately injected penicillin or 1.5 hour delayed injected penicillin in protecting against β-hemolytic Streptococci:

Penicillin Level	DELAY		PONSE Cured
1/8	None	0	6
	1.5 hr	0	8
1/4	None	3	3
	1.5 hr	1	6
1/2	None	6	0
	1.5 hr	2	4
1	None	5	1
	1.5 hr	6	0
4	None	2	0
	1.5 hr	5	0

- Note, there are many 0 cells in the table; may have problems with the large sample normal approximations.
- We want to estimate the (common) OR between Delay and Response, given strata (Penicillin).
- We can consider the data as arising from J = 5,  $(2 \times 2)$  tables, where J = 5 penicillin levels.
- In the  $k^{th}$  row of  $(2 \times 2)$  table *j*, we assume the data have the following logistic regression model:

 $\begin{array}{l} \text{logit} \left( P[\text{Cured}|\text{penicillin}=j,\text{DELAY}_k] \right) = \\ \mu + \alpha_j + \beta \text{DELAY}_k, \end{array}$ 

where

$$\mathsf{DELAY}_k = \begin{cases} 0 \text{ if none} \\ 1 \text{ if } 1.5 \text{ hours} \end{cases}$$

• There are problems with the unconditional (usual) MLE, as we'll see in the computer output.

• This is equivalent to the logistic model

```
\begin{split} & \mathsf{logit}\{P[\mathsf{Cured}|\mathsf{penicillin}=j,\mathsf{DELAY}]\} = \\ & \mu + \alpha_1 z_1 + \ldots + \alpha_4 z_4 + \beta \mathsf{DELAY}_k, \end{split}
```

where

$$z_j = \begin{cases} 1 \text{ if penicillin} = j \\ 0 \text{ if otherwise} \end{cases}$$

- Note, we have not taken the ordinal nature of penicillin into account
- We are considering penicillin (and the associated dose) as a nuisance parameter and are not interested in drawing inference about penicillin
- We are interested in learning about the timing to give any dose of penicillin

# SAS Proc Logistic

• Using SAS Proc Logistic, we get the following results:

```
data cmh;
 input pen delay response count;
cards;
0.125 0 0 0 /* response: 1 = cured, 0 = died */
0.125 0 1 6 /* delay: 0 = none, 1 = 1.5 hrs */
0.125 1 0 0
0.125 1 1 8
0.250 0 0 3
0.250 0 1 3
0.250 1 0 1
0.250 1 1 6
0.500 0 0 6
0.500 0 1 0
0.500 1 0 2
0.500 1 1 4
1.000 0 0 5
1.000 0 1 1
1.000 1 0 6
1.000 1 1 0
4.000 0 0 2
4.000 0 1 0
4.000 1 0 5
4.000 1 1 0
```

```
proc logistic descending;
class pen(PARAM=ref);
  model response = pen delay ;
  freq count;
run;
```

/\* SELECTED OUTPUT \*/

Model Convergence Status

Quasi-complete separation of data points detected.

WARNING: The maximum likelihood estimate may not exist. WARNING: The LOGISTIC procedure continues in spite of the above warning. Results shown are based on the last maximum likelihood iteration. Validity the model fit is questionable.

Paramet	er	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq
Interce	pt	1	-13.7143	163.8	0.0070	0.9333
pen	0.125	1	25.1650	199.4	0.0159	0.8996
pen	0.25	1	13.6947	163.8	0.0070	0.9334
pen	0.5	1	11.9500	163.8	0.0053	0.9419
pen	1	1	10.0640	163.8	0.0038	0.9510
delay		1	1.8461	0.9288	3.9508	0.0468

WARNING: The validity of the model fit is questionable.

Odds Ratio Estimates

				Point	95% V	Vald
Effect		Estimate	Confidenc	ce Limits		
pen	0.125	vs	4	>999.999	<0.001	>999.999
pen	0.25	vs	4	>999.999	<0.001	>999.999
pen	0.5	vs	4	>999.999	<0.001	>999.999
pen	1	vs	4	>999.999	<0.001	>999.999
delay				6.335	1.026	39.115

#### Possible remedies

There are (at least) two possible remedies to this problem:

- 1. Assign a score  $w_j$  to penicillin level j and use this with logistic regression, or
- 2. Since we are interested in

$$\beta = \log(OR^{\mathsf{RESPONSE},\mathsf{DELAY}}),$$

an alternative is to eliminate the nuisance parameters (the effects of PENICILLIN LEVEL) by using conditional logistic regression.

Revised data step:

```
data cmh;
/* delay = 0 if None, 1 if 1.5 hr */
/* response = 1 if Cured, 0 if Died */
 input pen delay response count;
 if pen = 0.125 then pen1 = 1; else pen1=0;
 if pen = 0.250 then pen2 = 1; else pen2=0;
 if pen = 0.500 then pen3 = 1; else pen3=0;
 if pen = 1.000 then pen4 = 1; else pen4=0;
pen cont = pen;
 cards;
 ... (same as before)
proc logistic descending;
model response = pen cont delay /aggregate scale=1;
 freq count;
run;
```

	Analy	sis of Maxim	num Likeliho	od Estimates	
	22		Standard	Wald	
Parameter	DF	Estimate	Error	Chi-Square	Pr > ChiSq
Intercept	1	2.2328	0.8051	7.6910	0.0055
pen_cont delay	1 1	-6.8591 1.6984	1.9781 0.8623	12.0238 3.8792	0.0005 0.0489
ueray	<b>–</b>	1.0904	0.0023	3.0792	0.0409
Odds Ratio Estimates					
	Poi	nt	95% Wald		
Effect	Estima	te Conf	idence Limi	ts	
pen_cont	0.0	01 <0.0	001 0.	051	
delay	5.4	65 1.0	08 29.	620	
Note: This model converged					

## Approach 2 - Elimination of Nuisance Parameters

For the prospective study we have, the rows (y<sub>j1+</sub> and y<sub>j2+</sub>) of each (2 × 2) table are fixed, and we have two independent binomials:
1) Row 1:

$$Y_{j11} \sim Bin(y_{j1+}, p_{j1}),$$

where

$$p_{j1} = P[\text{Cured}|\text{penicillin}=j, \text{DELAY}=1]$$

and 2) Row 2:

$$Y_{j21} \sim Bin(y_{j2+}, p_{j2}),$$

where

$$p_{j2} = P[\text{Cured}|\text{penicillin}=j, \text{DELAY}=0]$$

• In terms of the logistic model,

$$\begin{array}{l} \text{logit} \left( P[\text{Cured} | \text{penicillin} = j, \text{DELAY} = 1] \right) = \\ \text{logit}(p_{j1}) = \mu + \alpha_j + \beta, \end{array}$$

and

$$\begin{array}{l} \text{logit}\left(P[\text{Cured}|\text{penicillin}=j,\text{DELAY}=0]\right) = \\ \text{logit}(p_{j2}) = \mu + \alpha_j \end{array}$$

• Then the log-odds ratio for the  $j^{th}$ ,  $(2 \times 2)$  table is

$$\log \frac{p_{j1}/(1-p_{j1})}{p_{j2}/(1-p_{j2})} = \log (p_{j1}) - \log (p_{j2})$$
$$= [\mu + \alpha_j + \beta] - [\mu + \alpha_j]$$
$$= \beta,$$

and, the odds ratio is

$$\frac{(OR^{\mathsf{RESPONSE},\mathsf{DELAY}}) = \frac{p_{j1}/(1-p_{j1})}{p_{j2}/(1-p_{j2})} = e^{\beta}}{p_{j2}/(1-p_{j2})}$$

## **Conditional Logistic Regression**

• The conditional likelihood is

$$L^{c}(\vec{y}) = \frac{\prod_{j=1}^{J} \binom{n_{j}}{y_{j}}}{\sum_{\vec{y} \in \Gamma} \prod_{j=1}^{J} \binom{n_{j}}{y_{j}}}$$

which is the product of the probability functions over the J tables or strata.

- We can find the conditional MLE,  $\hat{\boldsymbol{\beta}}_{CMLE}$  by maximizing  $L^{c}(\beta)$ .
- As with a single table, in 'large' samples,  $\hat{\beta}$  will be approximately unbiased, approximately normal, and have variance equal to the negative inverse of the (expected) second derivative.

For proper properties of the conditional likelihood  $L^{c}(\beta)$ , we need 'large' samples.

There are two kinds of 'large samples' that allow  $\hat{\beta}_{CMLE}$  to be approximately normal (by the central limit theorem).

- 1. We can have the within stata sample size  $(y_{j++})$  large (with the number of strata *J* being small).
- 2. As the number of strata, *J*, becomes large (we can  $y_{j++} < \infty$ , and, actually, we can have  $y_{j++}$  as small as 2, which we will discuss later).

As a result, we can have both  $(y_{j++})$  large and J large, and this would be a large sample as well.

 This is in contrast to the usual (unconditional MLE); if you look at the logistic regression model

 $logit{P[Died|penici=j, DELAY_k]} =$ 

 $\mu + \alpha_j + \beta \mathsf{DELAY}_k,$ 

in order to be able to estimate each  $\alpha_j$ , we need each  $y_{j++}$  to be large

• In particular, if J is large, using unconditional logistic regression for

$$\operatorname{logit}(p_{jk}) = \mu + \alpha_j + \beta x_k,$$

we will have a large number of parameters ( $\alpha_j$ 's) to estimate, and we cannot unless each  $y_{j++}$  is also large.

• If  $y_{j++}$  is also large, then

$$\widehat{\beta}_{CMLE} \approx \widehat{\beta}_{MLE},$$

but the MLE is preferable since it is simpler computationally, and has nice properties.

## Mantel-Haenzsel Estimate

- One other possibility is the Mantel Haenszel estimate of a common odds ratio for a set of J,  $(2 \times 2)$  tables,
- As with the conditional MLE, the Mantel-Haenszel estimator of the common odds ratio is (asymptotically unbiased) if either  $(y_{j++})$  is large or J is large or both.
- However, the Mantel-Haenzel estimate does not generalize to data with many covariates, whereas the conditional likelihood does.

However, for the data studied today, we can calculate the CMH (or just the Mantel-Haenzsel)

```
proc freq;
tables pen*delay*response/cmh;
weight count;
run;
```

Estimates of the Common Relative Risk (Row1/Row2)

Type of Study	Method	Value	95% Confide	nce Limits
Case-Control	Mantel-Haenszel	4.6316	0.9178	23.3731
(Odds Ratio)	Logit **	3.9175	0.6733	22.7925

\*\* These logit estimators use a correction of 0.5 in every cell of those tables that contain a zero. Tables with a zero row or a zero column are not included in computing the logit estimators.

# **Conditional Logistic Regression**

- We can use SAS Proc Logistic as before to do conditional logistic regression.
- SAS Proc Logistic will give us the conditional logistic regression estimate of the odds ratio, and an exact 95% confidence interval for the odds ratio using the conditional likelihood.

```
proc logistic descending;
class pen(PARAM=ref);
  model response = pen delay ;
  exact delay / estimate = both /*both = logor & or */;
  freq count;
run;
```

```
/* SOME OUTPUT */
```

Conditional Exact Tests

			p-Va	alue
Effect	Test	Statistic	Exact	Mid
delay	Score	4.0880	0.0587	0.0384
_	Probability	0.0407	0.0587	0.0384

Exact Parameter Estimates						
Parameter	95% Confidence Estimate Limits p-Value					
delay	1.6903	-0.2120	4.1588	0.0937		
Exact Odds Ratios						
Parameter	Estimate	95% Confi Limit		p-Value		
delay	5.421	0.809	63.995	0.0937		

# Effectiveness of immediately injected penicillin or 1.5 hour

#### OR of Delayed injected penicillin in protecting against Streptococci

Variable	ODDS Ratio	95% Confidence Lower	Limits Upper
COND. LOGIT	 5.421	0.809	63.995
MANTEL-HA	3.918	0.673	22.793
LOGISTIC*	6.335	1.026	39.115
LOGISTIC**	5.465	1.008	29.620

- From each, we see that odds of being cured if there is no delay is about 5 times of what it is if there is a delay.
- Although we get an estimate from (ordinarly) logistic regression, it is not as stable as the other methods.
- Note, these strata sizes were small, but there are situations when they are even smaller (matched pairs, discussed later).
- First though, let's look at a couple of other models.
- In particular, let's test for interaction.

# Penicillin Data

- If we fit the model with interaction between delay and penicillin, (with 4 dummy variables for penicillin and the 4 for the interaction), and the unconditional had the same problems as above. The conditional also has problems:
- Proc Logistic gave me the following:

```
proc logistic descending;
class pen(PARAM=ref);
  model response = pen delay pen*delay ;
  exact pen*delay /estimate;
  freq count;
run;
```

delay\*pen 1

Exact Conditional Analysis

Conditional Exact Tests --- p-Value ---Effect Statistic Mid Exact Test 7.0666 0.1354 0.0990 delay\*pen Score Probability 0.0729 0.1354 0.0990 Exact Parameter Estimates 95% Confidence Estimate Limits Parameter p-Value delay\*pen 0.125 .# .# delay\*pen 0.25 • • delay\*pen 0.5 .# •

NOTE: # indicates that the conditional distribution is degenerate.

.#

- Proc Logistic can give us an exact conditional Score test (with p-value = .1354), even though it cannot give us estimates of the interaction terms.
- Thus, I fit a little simpler model, with the main effects of penicillin fitted using dummy variables, but the interaction with penicillin continuous  $(w_j)$ :

 $\begin{array}{l} \text{logit} \left( P[\text{Cured} | \text{penici} = j, \text{DELAY} = k] \right) = \\ \mu + \alpha_j + \beta_1 x_k + \beta_{12} x_k w_j, \end{array}$ 

where

$$x_k = \begin{cases} 0 \text{ if DELAY} = \text{none } (k = 1) \\ 1 \text{ if DELAY} = 1.5 \text{ hours } (k = 2) \end{cases}$$

and  $w_j$  equals the actual dose (0.125, .25, .5, or 1), treated as continuous.

• We are interested in testing for homogeneity of the odds ratio across the J stratum.

```
proc logistic descending;
class pen(PARAM=ref);
  model response = pen delay pen_cont*delay ;
  exact pen_cont*delay /estimate;
  freq count;
run;
```



		95% Conf	idence	
Parameter	Estimate Limits		p-Value	
delay*pen_cont	-4.3902	-15.4187	2.8438	0.2708

• When W (pennicilin) is treated as continuous,

```
\begin{array}{l} \text{logit} \left( P[\text{Cured} | \text{penici} = j, \text{DELAY} = k] \right) = \\ \mu + \gamma_1 w_j + \beta_1 x_k + \beta_{12} x_k w_j \end{array}
```

the unconditional logistic regression also converged, and we get the following comparison

95% Confide	nce Limits		
PEN	*DELAY		
METHOD	Parameter Estimate	P Value	
COND	-4.3902	0.2708	
UNCOND	-3.4990	0.4288	

- Although both algorithms converged, the estimates are pretty large, and the confidence intervals are very wide (not shown).
- However, we have indication that the common odds ratio assumption is divalid gistic Regression p. 39/40

### LogXact

- All of the models could have been fit in LogXact, a software package specializing in conditional (or exact) logistic regression
- I noted small variations in the parameter estimates between SAS and LogXact.
- LogXact is mouse driven and is easy to use
- SAS is easier to do model building activities creation of derived variables, "selection ="routines and "copy and paste" model code