## STAT 371: LAB MANUAL: WINTER 2015 MACEWAN UNIVERSITY KATHLEEN LAWRY-BATTY

#### **Contingency Tables**

A contingency table is a place where a joint distribution of two categorical variables, say X and Y, where X has I levels and Y has J levels, can be summarized. It is usual to consider an X variable to be an explanatory variable, and a Y variable a response variable. Both marginal distributions and conditional distributions can be of interest. Data that is used can either be offered in raw form (Format 1) or in summarized form (Format 2) in SPSS.

Example:

The file HockeyGenderRaw contains two variables. The response variable WatY (Yes-1, No-2) records whether or not an individual watches hockey and appears in column 1 and the explanatory variable GenX (Male -1, Female-2) records gender and appears in column 2.

(FORMAT 1)			
WatY	GenX		
1	1		
1	1		
1	2		
1	2		
1	2		
2	1		
2	1		
2	1		
2	2		
2	2		
2	2		

The file WatchHockeyGenderSummarized summarizes the counts of observations for the (GenX, WatY) pairs (1,1), (1,2), (2,1) and (2,2).

#### (FORMAT 2)

GenX	Wat	Y	Count
1	1	33	
1	2	11	
2	1	9	
2	2	36	

When working with raw data with columns of categorical variables, one can create a summary cross-tab table that counts the groupings. Here are the commands to do so for the hockey watching and gender raw data. <u>Commands</u>

Analyze>Descriptive Statistics>Crosstabs Row(s): GenX Column(s): WatY OK

#### GenX \* WatY Crosstabulation

			Wa		
			1	2	Total
GenX	1	Count	33	11	44
		% within GenX	75.0%	25.0%	100.0%
	2	Count	9	36	45
		% within GenX	20.0%	80.0%	100.0%
Total		Count	42	47	89
		% within GenX	47.2%	52.8%	100.0%

#### 3 ways of comparing two categorical variables (each with 2 levels)

It is usual to have X in the row, and Y in the column in a contingency table. Here  $n_{ij}$  represents the count of observations in the ijth cell, i=1,2 and j=1,2, while  $n_{i+}$  is the count in row I,  $n_{+J}$  is the count in column j, and n is the total number of observations.

		Y(Respo	onse)	
		1	2	Total
X(Explanatory)	1	<b>n</b> <sub>11</sub>	<b>n</b> <sub>12</sub>	<b>n</b> <sub>1+</sub>
	2	<b>n</b> 21	n <sub>22</sub>	<b>n</b> 2+
	Total	n <sub>+1</sub>	n <sub>2+</sub>	n

We write:

$$\begin{split} P(Y=1|X=1) &= p_1 = n_{11}/n_{1+} \\ P(Y=1|X=2) &= p_2 = n_{21}/n_{2+} \end{split}$$

**Compare proportions**:

The Wald (1- $\alpha$ )% CI is:  $(p_1-p_2) \pm z_{\alpha/2} \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$ 

#### **Relative Risk**:

 $\mathbf{rr} = \pi_1 / \pi_2, \ \widehat{\mathbf{rr}} = p_1 / p_2$ 

The distribution of  $\hat{rr}$  is very skewed, so we look at the distribution of  $\ln(\hat{rr})$  instead.

The CI for ln(rr) is :  $\ln(p_1/p_2) \pm z_{\alpha/2} \sqrt{\frac{1-p_1}{n_1p_1} + \frac{1-p_2}{n_2p_2}}$ 

Convert (exponentiate) the endpoints to find a CI for rr.

#### **Odds**

Odds<sub>i</sub>(of success in row i)=  $p_i/(1-p_i)$ , Hence:  $p_i = odds_i/(odds_i + 1)$ Odds Ratio:  $\theta = odds_1/odd_2$   $\hat{\theta} = n_{11}n_{22}/n_{12}n_{21}$  (both variables response variables)

 $\hat{\boldsymbol{\theta}}$  is MLE of  $\boldsymbol{\theta}$ 

Distribution of  $\widehat{\boldsymbol{\theta}}$  very skewed, so look at distribution of  $\ln(\widehat{\boldsymbol{\theta}})$  instead

CI: for  $\ln(\theta)$ :  $\ln(\hat{\theta}) \pm z_{\alpha/2} \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{22}} + \frac{1}{n_{12}} + \frac{1}{n_{21}}}$ 

Convert (exponentiate) the endpoints to find a CI for  $\theta$ .

#### Class Example:

In SPSS		
PCR	Relapse	Total
1	1	30
1	2	45
2	1	8
2	2	95
	1	

PCR is X (explanatory variable)

Relapse is Y (response variable)

#### **SPSS** Commands (to produce contingency table)

Data>Weight Cases \* Weight by Frequency Variable: Total OK Analysis>Descriptive Statistics>Crosstabs Row(s): PCR Column(s): Relapse Statistics popup: ✓ Risk Cells popup: Counts: ✓ Observed Percentages: ✓ Row

\*Note: If data is in two columns, one for PCR (X) and one for Relapse (Y), the weight cases command lines are unnecessary.

PCR * Relapse Crosstabulation					
			Rela	pse	
			1	2	Total
PCR	1	Count	30	45	75
		% within PCR	40.0%	60.0%	100.0%
	2	Count	8	95	103
		% within PCR	7.8%	92.2%	100.0%
Total		Count	38	140	178
		% within PCR	21.3%	78.7%	100.0%

Risk Estimate

		95% Confidence Interva	
	Value	Lower	Upper
Odds Ratio for PCR (1 / 2)	7.917	3.361	18.648
For cohort Relapse = 1	5.150	2.504	10.590
For cohort Relapse = 2	.651	.536	.789
N of Valid Cases	178		

$$P(Y=1|X=1) = p_1 = n_{11}/n_{1+} = 0.4$$
,  $P(Y=1|X=2) = p_2 = n_{21}/n_{2+} = 0.078$ 

(1-
$$\alpha$$
)% Wald CI for (p<sub>1</sub>-p<sub>2</sub>) = (p<sub>1</sub>-p<sub>2</sub>) ±  $z_{\alpha/2} \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$ 

95% CI is (.200, 0.445) (this must be calculated by hand) Note that 0 is not in the 95% CI.

With 95% confidence, we have sufficient evidence that the proportion of PCR positive children who relapse is between 0.2 and 0.445 larger than the proportion of PCR- children who relapse.

Another way of saying this:

With 95% confidence, we have sufficient evidence that the proportion of children who relapse if they have a positive PCR is between 0.2 and 0.445 larger than the proportion of children who relapse if they have a negative PCR.

#### Estimated relative risk = $\widehat{\mathbf{rr}} = \mathbf{p_1/p_2} = 0.4/0.78 = 5.15$ ,

95% CI for rr is (2.504, 10.590) from above. (exponentiated endpoints of CI for  $ln(p_1/p_2)$ )

Note that 1 is not in the 95% CI.

With 95% confidence, we have sufficient evidence that probability of a relapse is between 2.5 and 10.6 times higher for children with a positive PCR than a negative PCR.

#### Estimated odds ratio = $\hat{\theta} = n_{11}n_{22}/n_{12}n_{21} = (30x95)/(8x45) = 7.917$

95% CI for  $\theta$  is (3.361, 18.648), from above. (exponentiated endpoints of CI for  $\ln(\hat{\theta})$ )

Note that 1 is not in the 95% CI.

With 95% confidence, we have sufficient evidence that a relapse is associated with PCR. The odds of a relapse is 3.361 to 18.648 X higher for the PCR+ group than for the PCR- group.

#### Odds Ratio (another option for commands for a 2x2 contingency table with 2 levels per category)

SPSS Commands Analysis>Descriptive Statistics>Crosstabs Row(s): PCR Column(s): Relapse Statistics popup: Cochran's and Mantel-Haenszel statistics Test common odds ratio equals: 1 Cells popup: Counts: ✓ Observed

<u>\*These commands can only produce a 95% confidence interval.</u> <u>\*These commands also produce a confidence interval for  $ln(\hat{\theta})$ </u>

#### Mantel-Haenszel Common Odds Ratio Estimate

Estimate			7.917
In(Estimate)			2.069
Std. Error of In(Estimate)			.437
Asymp. Sig. (2-sided)			.000
Asymp. 95% Confidence	Common Odds Ratio	Lower Bound	3.361
Interval		Upper Bound	18.648
	In(Common Odds Ratio)	Lower Bound	1.212
		Upper Bound	2.926

The Mantel-Haenszel common odds ratio estimate is asymptotically normally distributed under the common odds ratio of 1.000 assumption. So is the natural log of the estimate.

#### Example: Gender/Hockey Data

Illustrate the gender and hockey watching variables in the three ways introduced above. ALWAYS BEST TO SET UP YOUR PROBLEM SO THAT  $p_1 = P(Y=1|X=1)$  (in upper left corner of table) is the largest probability you expect to find.

X (explanatory) = Gender, Y(response) = Watch Hockey

$$P(Y=1|X=1) = p_1 = n_{11}/n_{1+} = 0.75, P(Y=1|X=2) = p_2 = n_{21}/n_{2+} = 0.2$$
  
Wald (1-\alpha)% CI for (p\_1-p\_2) = (p\_1-p\_2) \pm z\_{\alpha/2} \sqrt{\frac{p\_1(1-p\_1)}{n\_1} + \frac{p\_2(1-p\_2)}{n\_2}} \rightarrow 95% CI becomes (...38, ...72)

0 is not in 95% CI. With 95% confidence, we have sufficient evidence that the proportion of males who watch hockey differs from the proportion of females who watch hockey. The proportion of males who watch hockey is between 38% and 72% higher than the proportion of females who watch.

Estimated relative risk =  $\hat{rr} = p_1/p_2 = 0.75/0.2 = 3.75$ 

95% CI for rr is (2.040, 6.893) ( exponentiated endpoints of CI for ln(p1/p2) )

1 is not in 95% CI. With 95% confidence, we have sufficient evidence that probability of being a hockey watcher is 2.040 to 6.893 times higher for males than females.

Estimated odds ratio =  $\hat{\theta} = n_{11}n_{22}/n_{12}n_{21} = (33X36)/(9x11) = 12$ 

95% CI for  $\theta$  is (4.416, 32.606)(exponentiated endpoints of CI for  $\ln(\hat{\theta})$ )

1 is not in 95% CI. With 95% confidence, sufficient evidence that a gender is associated with regular hockey watching.

We can also say: The odds of being a hockey watcher are 4.416 to 32.606 times higher for a male than a female.

#### GenX \* WatY Crosstabulation

			Wa		
			1	2	Total
GenX	1	Count	33	11	44
		% within GenX	75.0%	25.0%	100.0%
	2	Count	9	36	45
		% within GenX	20.0%	80.0%	100.0%
Total		Count	42	47	89
		% within GenX	47.2%	52.8%	100.0%

Risk Estimate

		95% Confidence Interva	
	Value	Lower	Upper
Odds Ratio for GenX (1 / 2)	12.000	4.416	32.606
For cohort WatY = 1	3.750	2.040	6.893
For cohort WatY = 2	.313	.184	.532
N of Valid Cases	89		

#### **Example: University Student Data: Coding**

Gen (F-Female, M-Male) Hsleep (hours sleep weeknight) Energy (scale of 1 to 10) HrHwkCrs (weekly) MTxH (minutes text per hour) Nsiblings (number of siblings) NChPlan (number of children have/planned) Nlang (number of spoken languages) Age (in years) HH3 (homework: low 0.5-1.5 $\rightarrow$ 1, med 2-3 $\rightarrow$ 2, high 3,5-5 $\rightarrow$ 3) EN3 (energy: low 2-4 $\rightarrow$ 1, med5-7 $\rightarrow$ 2, high 8+ $\rightarrow$ 3) HH2 (homework: low 0-1.5 $\rightarrow$ 1, high 2+ $\rightarrow$ 2) HS3 (sleep, low 4-5 $\rightarrow$ 1, med 6-7 $\rightarrow$ 2, high 8-9 $\rightarrow$ 3) HS2 (sleep  $\rightarrow$ , low 4-6  $\rightarrow$ 1, high 7-9 $\rightarrow$ 2) NL2 (languages: unilingual -1, multilingual -2) University Student Data (how to write up findings if you do not have significance Is gender related to number of languages spoken? Gender (rows) (1 - Female, 2 - Male) NL2: Languages (columns)(1 - Unilingual, 2 - Multilingual)

Page .

#### Gender \* NL2 Crosstabulation

			NI	_2	
			1.00	2.00	Total
Gender	F	Count	38	15	53
		% within Gender	71.7%	28.3%	100.0%
	М	Count	35	18	53
		% within Gender	66.0%	34.0%	100.0%
Total		Count	73	33	106
		% within Gender	68.9%	31.1%	100.0%

Risk Estimate

		95% Confidence Interval		
	Value	Lower	Upper	
Odds Ratio for Gender (F / M)	1.303	.571	2.972	
For cohort NL2 = 1.00	1.086	.840	1.403	
For cohort NL2 = 2.00	.833	.471	1.473	
N of Valid Cases	106			

X (explanatory) = Gender (1 (Female),2 (Male)) Y(response) = NL2(1 (unilingual), 2 (multilingual)) P(Y=1|X=1) =  $p_1 = n_{11}/n_{1+} = 0.72$ , P(Y=1|X=2) =  $p_2 = n_{21}/n_{2+} = 0.66$ <u>Wald (1- $\alpha$ )% CI for ( $p_1$ - $p_2$ ) = ( $p_1$ - $p_2$ )  $\pm z_{\alpha/2} \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} \rightarrow 95\%$  CI becomes (-.13,.25)</u>

0 is in 95% CI. For  $\alpha = 0.05$ , we do not have sufficient evidence that the proportion of females who are unilingual differs from the proportion of males who are unilingual. The proportion of females who are unilingual is between 13% less to 25% higher than the portion of males who are unilingual.

Estimated relative risk =  $\hat{rr} = p_1/p_2 = 0.72/0.66 = 1.086$ 

95% CI for rr is (0.840,1.403) (exponentiated endpoints of CI for ln(p1/p2))

1 is in 95% CI. For  $\alpha = 0.05$ , we do not have significant evidence that the probability of being unilingual is a significant number of times higher for females than males. The probability of being unilingual is between 1.190 (1/.840) times higher for a male to 1.403 times higher for a female.

Estimated odds ratio =  $\hat{\theta} = n_{11}n_{22}/n_{12}n_{21} = (38X18)/(35x15) = 1.303$ 

95% CI for  $\theta$  is (0.571, 2.972) (exponentiated endpoints of CI for  $\ln(\hat{\theta})$ )

1 is in 95% CI. For  $\alpha = 0.05$ , we do not have sufficient evidence that gender is associated with language plurality. The odds of being unilingual is between 1.751 (1/.571) times higher for a male than a female to 2.972 times higher for a female than a male.

#### Possible Useful SPSS Commands for transforming/editing data

Use Transform > Recode into Different Variables to get new recoded explanatory variable Use Data > Select Cases to get subset worksheets Use copy/paste to create column for new explanatory variable

#### **Data: Class Examples Flu/Vaccine**

EXAMPLE 17 and 18 DATA	No Vaccine (1)	One Shot (2)	Two Shots (3)	Total
Flu (1)	24	9	13	46
No Flu (2)	289	100	565	954
Total	313	109	578	1000

Input Data to SPSS: Flu (Response), Vaccine (Explanatory):

Here data is presented counter to the usual way, with the response variable on the rows and the Explanatory variable on the columns. Data is in fluvaccine.sav

<u>Flu</u>	Vaccine	Count
1	1	24
1	2	9
1	3	13
2	1	289
2	2	100
2	3	565

We illustrate how to create a cross-tab table for flu and vaccine, and perform a test of independence. Standardized residuals are calculated to identify cells of possible influence.

 SPSS Commands

 Data>Weight Cases

 Weight by Frequency Variable: Count

 OK

 Analysis>Descriptive Statistics>Crosstabs

 Row(s): Flu

 Column(s): Vaccine

 Statistics popup:
 ✓ Chi-square

 Cells popup: Counts:
 ✓ Observed

 ✓ Expected
 ✓ Residuals: Adjusted standardized\*

\*Adjusted Standardized Residuals in SPSS are Standardized Cell Residuals in Class Notes

#### **Output: Class Example Flu/Vaccine**

				Vacine			
			1	2	3	Total	
Flu	1	Count	24	9	13	46	
		Expected Count	14.4	5.0	26.6	46.0	
		Adjusted Residual	3.1	1.9	-4.2		
	2	Count	289	100	565	954	
		Expected Count	298.6	104.0	551.4	954.0	
		Adjusted Residual	-3.1	-1.9	4.2		
Total		Count	313	109	578	1000	
		Expected Count	313.0	109.0	578.0	1000.0	

Flu \* Vacine Crosstabulation

#### Chi-Square Tests

	Value	df	Asymp. Sig. (2-sided)
Pearson Chi-Square	17.313 <sup>a</sup>	2	.000
Likelihood Ratio	17.252	2	.000
Linear-by-Linear Association	14.915	1	.000
N of Valid Cases	1000		

a. 0 cells (.0%) have expected count less than 5. The minimum expected count is 5.01.

#### **Test of Independence**

#### Question: Is catching flu and vaccination status associated?

Ho: Flu and Vaccine independent vs Ha: Ho not true Ho: Flu and Vaccine not associated vs Ha:Ho not true  $\alpha = 0.05$ Assume: random sample, all  $\mu_{ij} >=5$ Test statistic: (HIGHLIGHTED ON OUTPUT) Pearson: 17.313 Likelihood Ratio: 17.252 Decision: Reject Ho (p-values are both <=.001) There is significant evidence that flu and vaccination status are associated.

Adjusted standardized residuals >= 1.96 in absolute magnitude identify cells contributing to the significance. SEE CELLS 11, 13, 31, AND 33 IN THE OUTPUT. Here people without the vaccine are more likely to get flu than people with 2 shots of the vaccine.

#### Tests of Independence: return to University Data

Question: Is hours of homework associated with energy level? Choose your coding and make new SPSS columns Column HH2: Low( $0-1.5 \rightarrow 1$ ), High( $2+ \rightarrow 2$ ) Column EN3: Low ( $2-4 \rightarrow 1$ ), Medium ( $5-7 \rightarrow 2$ ), High ( $8+ \rightarrow 3$ ) Perform a test of independence in SPSS and write up your hypothesis test.

#### SPSS Commands

Analysis>Descriptive Statistics>Crosstabs

Row(s): HH2

Column(s): EN3

Statistics popup: ✓ Chi-square

Cells popup: Counts: ✓ Observed

✓ Expected

✓ Residuals: Adjusted standardized

### HH2 \* EN3 Crosstabulation

				EN3			
			1.00	2.00	3.00	Total	
HH2	1.00	Count	8	23	16	47	
		Expected Count	14.2	21.3	11.5	47.0	
		% within HH2	17.0%	48.9%	34.0%	100.0%	
		Adjusted Residual	-2.6	.7	2.0		
	2.00	Count	24	25	10	59	
		Expected Count	17.8	26.7	14.5	59.0	
		% within HH2	40.7%	42.4%	16.9%	100.0%	
		Adjusted Residual	2.6	7	-2.0		
Total		Count	32	48	26	106	
		Expected Count	32.0	48.0	26.0	106.0	
		% within HH2	30.2%	45.3%	24.5%	100.0%	

#### Chi-Square Tests

	Value	df	Asymp. Sig. (2-sided)
Pearson Chi-Square	8.215 <sup>a</sup>	2	.016
Likelihood Ratio	8.491	2	.014
Linear-by-Linear Association	7.911	1	.005
N of Valid Cases	106		

a. 0 cells (.0%) have expected count less than 5. The minimum expected count is 11.53.

#### University Data: Test of Independence for Hours Homework and Energy Level

Ho: Homework and Energy independent vs Ha: Ho not true

Ho: Homework and Energy not associated vs Ha:Ho not true

#### $\alpha = 0.05$

Assume: random sample, all  $\mu_{ij} >= 5$ 

Test statistic: Pearson  $\chi^2$ : 8.215 Likelihood Ratio: 8.491 Decision: Reject Ho (p-values are both <=.05)

There is significant evidence that Homework and Energy are associated. The adjusted standardized residuals >= 1.96 in absolute magnitude identify cells contributing to the significance. SEE CELLS 11, 13, 21, AND 23 IN THE OUTPUT. Here people doing less homework are more

likely to have more energy and people doing more homework are more likely to be less energetic.

#### **Test of Linear Trend: Ordinal Variables (looking for linear trend)** Data: Class Example on tartar and smoking levels

By hand	Tartar			
Smoking	None	Middle	Heavy	Total
No	284	236	48	568
Middle	606	983	209	1798
Heavy	1028	1871	425	3324
Total	1918	3090	682	5690

#### IN SPSS (Smoketartar.sav)

Smoke	Tartar	Count	Smoke2	Smoke3
1.00	1.00	284.00	.0	284.50
1.00	2.00	236.00	.0	284.50
1.00	3.00	48.00	.0	284.50
2.00	1.00	606.00	2.00	1467.50
2.00	2.00	983.00	2.00	1467.50
2.00	3.00	209.00	2.00	1467.50
3.00	1.00	1028.00	3.00	4028.50
3.00	2.00	1871.00	3.00	4028.50
3.00	3.00	425.00	3.00	4028.50

Here Smoking (X) on the rows explains Tartar (Y) on the columns Tarter has 3 levels (none (1), middle (2), heavy (3) Various numerical scores can reflect distances between smoking levels (No, Middle, Heavy) Smoke: Scores are 1, 2, 3 for none, middle, heavy smoking Smoke 2: Scores are 0, 2, 3 for none, middle, heavy smoking Smoke 3: Scores are midranks for none, middle, heavy smoking

#### SPSS Commands

Data>Weight Cases Weight by Frequency Variable: Count OK Analysis>Descriptive Statistics>Crosstabs

#### Row(s): Smoke1, Smoke2, Smoke3 (3 POSSIBLE SCORES)

Column(s): Tartar

Statistics popup:

- ✓ Chi-square
- $\checkmark$  Correlation
- ✓ Kendall's Tau-b

Cells popup: Counts: ✓ Observed

Pearson's r changes depending on scores, but generally "gives similar results" (Agresti) Spearman's R takes the fact that the data is ordinal into account.

Kendal's  $\tau$ -b examines the relationship between concordant and discordant sample observations.

Spearman's R and Kendall's Tau-b do not change with different scores.

Smoke1	Chi-Square Tests			Symmetric Measures						
		Value	df	Asymp. Sig. (2-sided)			Value	Asymp. Std. Error <sup>a</sup>	Approx. T <sup>b</sup>	Approx. Sig.
	Pearson Chi-Square	79.717 <sup>a</sup>	4	.000	Ordinal by Ordinal	Kendall's tau-b	.080	.012	6.437	.000
	Likelihood Ratio	76 226	4	000		Spearman Correlation	.086	.013	6.545	°000.
	Lincor by Lincor	51 212		.000	Interval by Interval	Pearson's R	.095	.013	7.188	°000.
	Association	51.213	1 1	.000	N of Valid Cases		5690			
	N of Valid Cases	5690			a. Not assuming	g the null hypothesis.				
	a. 0 cells (.0%) have expected count is 68.08.	ected count l	ess than 5. 1	The minimum	b. Using the asy c. Based on nor	mptotic standard error as: mal approximation.	suming the	null hypothesis.		
Smoke2	Cł	ni-Square Te	sts			Symn	netric Meas	ures		
		Value	df	Asymp. Sig. (2-sided)			Value	Asymp. Std. Error <sup>a</sup>	Approx. T <sup>b</sup>	Approx. Sig.
	Pearson Chi-Square	79.717 <sup>a</sup>	4	.000	Ordinal by Ordinal	Kendall's tau-b	.080	.012	6.437	.000
	Likelihood Ratio	76.226	4	.000		Spearman Correlation	.086	.013	6.545	.000°
	Linear by Linear	60.930	1	000	Interval by Interval	Pearson's R	.103	.013	7.841	.000*
	Association	00.000	· ·		a Not assuming	the null hynothesis	5050			
	N of Valid Cases	5690			b. Using the asy	mptotic standard error ass	umina the n	ull hypothesis.		
	a. 0 cells (.0%) have exp expected count is 68.08.	ected count	less than 5.	The minimum	c. Based on norr	nal approximation.				
Smoke3	Ch	i-Square Tes	sts			Symm	etric Measu	ires		
		Value	df	Asymp. Sig. (2-sided)			Value	Asymp. Std. Error <sup>a</sup>	Approx. T <sup>b</sup>	Approx. Sig.
	Pearson Chi-Square	79.717 <sup>a</sup>	4	.000	Ordinal by Ordinal	Kendall's tau-b	.080	.012	6.437	.000
	Likelihood Ratio	76.226	4	000		Spearman Correlation	.086	.013	6.545	.000°
	Linear by Linear	20 710		000	Interval by Interval	Pearson's R	.084	.013	6.324	.000°
	Association	55.715	'	.000	a Not assuming	the null hypothesis	2090			
	N of Valid Cases	5690			b. Using the asyr	nptotic standard error assu	iming the nu	ll hypothesis.		
	a. 0 cells (.0%) have exp expected count is 68.08.	ected count l	ess than 5. 1	he minimum	<ul> <li>b. Using the asymptotic standard error assuming the null hypothesis.</li> <li>c. Based on normal approximation.</li> </ul>					

#### Ho: $\rho = 0$ versus Ha: $\rho \neq 0$ Random Samples

	Test statistic: $M_o^2 = (n-1)r^2$ , df = 1	$Pvalue = P(\chi^2 > M_o^2)$
Pearson's R	Smoke1: 0.095, Smoke2: 0.103, Smoke3: 0.084	< 0.001
Spearman	0.086	< 0.001
Kendall's τ – b	0.080	< 0.001

There is significant evidence that there is a linear trend association. The more a person smokes the more tartar builds up.

#### **Three Way Table:**Class Example (alienwallethreewaydata.sav)

This example considers 3 categorical variables, X, Y, and Z, where we can control Z by looking at partial tables for X and Y for each level of Z.

- Y response variable (Rating(score as 1,2,3,4,5)
- X explanatory variable (Movie: Alien (1) or Wall-e(2))
- Z-confounding variable (Gender: Female (1) and Male (2)

The data in the data file is:

Movie Freq	Rating	Sex	Score	Sex_n	Movie_n
Alien 430.00	12	f	1.00	1.00	1.00
Wall-e 1104.00	12	f	1.00	1.00	2.00
Alien 356.00	34	f	2.00	1.00	1.00
Wall-e 441.00	34	f	2.00	1.00	2.00
Alien 1313.00	56	f	3.00	1.00	1.00
Wall-e 1274.00	56	f	3.00	1.00	2.00
Alien 5075.00	78	f	4.00	1.00	1.00
Wall-e 5998.00	78	f	4.00	1.00	2.00
Alien 6777.00	910	f	5.00	1.00	1.00
Wall-e 19106.00	910	f	5.00	1.00	2.00
Alien 1502.00	12	m	1.00	2.00	1.00
Wall-e 5654.00	12	m	1.00	2.00	2.00
Alien 1404.00	34	m	2.00	2.00	1.00
Wall-e 2261.00	34	m	2.00	2.00	2.00
Alien 7090.00	56	m	3.00	2.00	1.00
Wall-e 8199.00	56	m	3.00	2.00	2.00
Alien 49158.00	78	m	4.00	2.00	1.00
Wall-e 43116.00	78	m	4.00	2.00	2.00
Alien 77097.00	910	m	5.00	2.00	1.00
Wall-e 94079.00	910	m	5.00	2.00	2.00

The tables are below

Male	Rating					
Movie	1-2	3-4	5-6	7-8	9-10	Total
Alien	1502	1404	7090	49158	77097	136251
Wall-e	5654	2261	8199	43116	94079	153309
Total	7156	3665	15289	92274	171176	289560
Female	Rating					

Movie	1-2	3-4	5-6	7-8	9-10	Total
Alien	430	356	1313	5075	6777	13951
Wall-e	1104	441	1274	5998	19106	27923
Total	1534	797	2587	11073	25883	41874
Both	Rating					
Movie	1-2	3-4	5-6	7-8	9-10	Total
Alien	1932	1760	8403	54233	83874	150202
Wall-e	6758	2702	9473	49114	113185	181232
Total	8690	4462	17876	103347	197059	331434

A  $\chi^2$  test for conditional independence (female and male) and correlation, r, to test for **conditional trend** (female and male) using the chosen scores 1, 2, 3, 4, 5 considers the data for each gender separately. SPSS commands to obtain useful output follow.

#### SPSS Commands

Data>Weight Cases Weight by Frequency Variable: Freq OK Analysis>Descriptive Statistics>Crosstabs Row(s): movie\_n Column(s): score Layer 1 of 1:sex\_n Statistics popup: ✓ Chi-square ✓ Correlation Cells popup: Counts: ✓ Observed

A  $\chi^2$  test for marginal independence and Correlation, r, to test **marginal trend** ((for chosen scores: 1,2,3,4,5) considers the data for both genders together. SPSS commands are as above, but no Layer(s) command is included.

#### Note the following:

- 1. Rows and Columns must be numerical in order for r to be calculated
- 2. Pearson's r will report the r corresponding to the scores chosen (here 1,2,3,4,5)
- 3. Spearman's r (based on mid-ranks scores)
- 4. Can calculate Chi-square statistic with Movie, Score, and Sex as string.

Output of interest is shown below. Necessary test statistic values and p-values are obtained from the output and summarized.

			movie_n	* score * se	ex Crosstab	ulation			
						score			
sex				1.00	2.00	3.00	4.00	5.00	Total
ſ	movie_n	1.00	Count	430	356	1313	5075	6777	13951
			% within movie_n	3.1%	2.6%	9.4%	36.4%	48.6%	100.0%
			% within score	28.0%	44.7%	50.8%	45.8%	26.2%	33.3%
		2.00	Count	1104	441	1274	5998	19106	27923
			% within movie_n	4.0%	1.6%	4.6%	21.5%	68.4%	100.0%
			% within score	72.0%	55.3%	49.2%	54.2%	73.8%	66.7%
	Total		Count	1534	797	2587	11073	25883	41874
			% within movie_n	3.7%	1.9%	6.2%	26.4%	61.8%	100.0%
			% within score	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%
m	movie_n	1.00	Count	1502	1404	7090	49158	77097	136251
			% within movie_n	1.1%	1.0%	5.2%	36.1%	56.6%	100.0%
			% within score	21.0%	38.3%	46.4%	53.3%	45.0%	47.1%
		2.00	Count	5654	2261	8199	43116	94079	153309
			% within movie_n	3.7%	1.5%	5.3%	28.1%	61.4%	100.0%
			% within score	79.0%	61.7%	53.6%	46.7%	55.0%	52.9%
	Total		Count	7156	3665	15289	92274	171176	289560
			% within movie_n	2.5%	1.3%	5.3%	31.9%	59.1%	100.0%
			% within score	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%

movie_n * score Crosstabulation										
					score					
			1.00	2.00	3.00	4.00	5.00	Total		
movie_n	1.00	Count	1932	1760	8403	54233	83874	150202		
		% within movie_n	1.3%	1.2%	5.6%	36.1%	55.8%	100.0%		
		% within score	22.2%	39.4%	47.0%	52.5%	42.6%	45.3%		
	2.00	Count	6758	2702	9473	49114	113185	181232		
		% within movie_n	3.7%	1.5%	5.2%	27.1%	62.5%	100.0%		
		% within score	77.8%	60.6%	53.0%	47.5%	57.4%	54.7%		
Total		Count	8690	4462	17876	103347	197059	331434		
		% within movie_n	2.6%	1.3%	5.4%	31.2%	59.5%	100.0%		
		% within score	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%		

movie_n * score Crosstabulation										
					score					
			1.00	2.00	3.00	4.00	5.00	Total		
movie_n	1.00	Count	1932	1760	8403	54233	83874	150202		
		% within movie_n	1.3%	1.2%	5.6%	36.1%	55.8%	100.0%		
		% within score	22.2%	39.4%	47.0%	52.5%	42.6%	45.3%		
	2.00	Count	6758	2702	9473	49114	113185	181232		
		% within movie_n	3.7%	1.5%	5.2%	27.1%	62.5%	100.0%		
		% within score	77.8%	60.6%	53.0%	47.5%	57.4%	54.7%		
Total		Count	8690	4462	17876	103347	197059	331434		
		% within movie_n	2.6%	1.3%	5.4%	31.2%	59.5%	100.0%		
		% within score	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%		

	Symmetric Measures										
		Value	Asymp. Std. Error <sup>a</sup>	Approx. T <sup>b</sup>	Approx. Sig.						
Interval by Interval	Pearson's R	006	.002	-3.267	.001°						
Ordinal by Ordinal	Spearman Correlation	.048	.002	27.571	°000.						
N of Valid Cases		331434									

a. Not assuming the null hypothesis.

b. Using the asymptotic standard error assuming the null hypothesis.

c. Based on normal approximation.

#### Chi-Square Tests

	Value	df	Asymp. Sig. (2-sided)
Pearson Chi-Square	4692.375 <sup>a</sup>	4	.000
Likelihood Ratio	4823.106	4	.000
Linear-by-Linear Association	10.671	1	.001
N of Valid Cases	331434		

a. 0 cells (.0%) have expected count less than 5. The minimum expected count is 2022.13.

sex		Value	df	Asymp. Sig. (2-sided)
ſ	Pearson Chi-Square	1793.102 <sup>a</sup>	4	.000
	Likelihood Ratio	1758.419	4	.000
	Linear-by-Linear Association	584.354	1	.000
	N of Valid Cases	41874		
m	Pearson Chi-Square	3778.475 <sup>b</sup>	4	.000
	Likelihood Ratio	3927.323	4	.000
	Linear-by-Linear Association	160.043	1	.000
	N of Valid Cases	289560		

**Chi-Square Tests** 

a. 0 cells (.0%) have expected count less than 5. The minimum expected count is 265.53.

b. 0 cells (.0%) have expected count less than 5. The minimum expected count is 1724.55.

#### Symmetric Measures

sex			Value	Asymp. Std. Error <sup>a</sup>	Approx. T <sup>b</sup>	Approx. Sig.
ſ	Interval by Interval	Pearson's R	.118	.005	24.344	°000.
	Ordinal by Ordinal	Spearman Correlation	.181	.005	37.564	.000°
	N of Valid Cases		41874			
m	Interval by Interval	Pearson's R	024	.002	-12.654	°000.
	Ordinal by Ordinal	Spearman Correlation	.029	.002	15.343	.000°
	N of Valid Cases		289560			

a. Not assuming the null hypothesis.

b. Using the asymptotic standard error assuming the null hypothesis.

c. Based on normal approximation.

# $_{\text{Page}}16$

<u>BY HAND</u> <u>SUMMARY</u>	$\chi^2$	Df	Р
Conditional Independence (female)	1793.102	4	<.0001
Conditional Independence (male)	3778.475	4	<0.001
Marginal independence	4692.375	4	<0.001

BY HAND SUMMARY	r	M <sup>2</sup>	df	Р
Trend (conditional female)	0.118	583.03	1	<0.001
Trend (conditional male)	-0.024	166.79	1	<0.001
Marginal trend	-0.006	11.93	1	0.001

#### Test of Independence: Write-Up

There is significant evidence that rating and movie are dependent (associated) in both the partial table for males and the partial table for females .....<u>so we say:</u>

#### Rating and Movie are conditionally dependent, given gender. OR

There is a conditional association of movie and rating for males. OR

#### There is a conditional association of movie and rating for females.

There is significant evidence that rating and movie are dependent when we look at the marginal table (that ignores male and female and puts all the counts together).

#### We have marginal dependence (association)of rating and movie.

#### Test of Linear Association: Write-Up

There is significant evidence of a difference in median for the Alien and Wall-e groups when we condition on females. The positive sign on r indicates that Wall-e has a higher median rating than Alien for the female group. There is significant evidence of a difference in medians for the Alien and Wall-e groups when we condition on males. The negative sign on r indicates that Wall-e has a lower median rating than Alien for the male group. There is significant evidence of a difference in medians for the Alien and Wall-e groups when we look at the entire group of males and females together (the marginal result). The negative sign on r indicates that Wall-e has a lower median rating sign on r indicates that Wall-e has a lower median rating sign on r indicates that Wall-e has a lower median rating sign on r indicates that Wall-e has a lower median rating sign on r indicates that Wall-e has a lower median rating sign on r indicates that Wall-e has a lower median rating sign on r indicates that Wall-e has a lower median rating sign on r indicates that Wall-e has a lower median rating sign on r indicates that Wall-e has a lower median rating sign on r indicates that Wall-e has a lower median rating than Alien for the marginal situation.

Tests of Independence and Linear Trend University Data

Question: Is hours of homework associated with energy level? Confounding Variable is Gender! Y: Response Variable: EN3: Low (2-4  $\rightarrow$ 1), Medium (5-7  $\rightarrow$ 2), High (8+  $\rightarrow$ 3) X: Explanatory Variable: HH2: Low( 0-1.5  $\rightarrow$  1), High(2+  $\rightarrow$  2) Z: Gender

#### (Female $\rightarrow 1$ , Male $\rightarrow 2$ )

#### Exercise: Perform a test of independence in SPSS. Write up your hypothesis test.

			HH2 * EN3 * Gende	er Crosstabu	lation		
					EN3		
Gende	er			1.00	2.00	3.00	Total
F	HH2	1.00	Count	6	7	5	18
			Expected Count	7.1	6.1	4.8	18.0
			Adjusted Residual	7	.5	.2	
		2.00	Count	15	11	9	35
			Expected Count	13.9	11.9	9.2	35.0
			Adjusted Residual	.7	5	2	
	Total		Count	21	18	14	53
			Expected Count	21.0	18.0	14.0	53.0
м	HH2	1.00	Count	2	16	11	29
			Expected Count	6.0	16.4	6.6	29.0
			Adjusted Residual	-2.7	2	2.9	
		2.00	Count	9	14	1	24
			Expected Count	5.0	13.6	5.4	24.0
			Adjusted Residual	2.7	.2	-2.9	
	Total		Count	11	30	12	53
			Expected Count	11.0	30.0	12.0	53.0

HH2 * EN3 Crosstabulation										
				EN3						
			1.00	2.00	3.00	Total				
HH2	1.00	Count	8	23	16	47				
		Expected Count	14.2	21.3	11.5	47.0				
		Adjusted Residual	-2.6	.7	2.0					
	2.00	Count	24	25	10	59				
		Expected Count	17.8	26.7	14.5	59.0				
		Adjusted Residual	2.6	7	-2.0					
Total		Count	32	48	26	106				
		Expected Count	32.0	48.0	26.0	106.0				

#### Chi-Square Tests

Geno	der	Value	df	Asymp. Sig. (2-sided)
F	Pearson Chi-Square	.486 <sup>a</sup>	2	.784
	Likelihood Ratio	.490	2	.783
	Linear-by-Linear Association	.244	1	.622
	N of Valid Cases	53		
М	Pearson Chi-Square	12.561 <sup>b</sup>	2	.002
	Likelihood Ratio	14.231	2	.001
	Linear-by-Linear Association	12.311	1	.000
	N of Valid Cases	53		

a. 1 cells (16.7%) have expected count less than 5. The minimum expected count is 4.75.

b. 1 cells (16.7%) have expected count less than 5. The minimum expected count is 4.98.

Symmetric Measures

Geno	ler		Value	Asymp. Std. Error <sup>a</sup>	Approx. T <sup>b</sup>	Approx. Sig.					
F	Ordinal by Ordinal	Kendall's tau-b	068	.128	533	.594					
		Spearman Correlation	072	.135	517	.608°					
	Interval by Interval	Pearson's R	068	.135	490	.626°					
	N of Valid Cases		53								
М	Ordinal by Ordinal	Kendall's tau-b	464	.094	-4.526	.000					
		Spearman Correlation	487	.099	-3.979	.000°					
	Interval by Interval	Pearson's R	487	.099	-3.977	.000°					
	N of Valid Cases		53								
a.	a Not assuming the null hypothesis										

Chi-Square	Tests

	Value	df	Asymp. Sig. (2-sided)
Pearson Chi-Square	8.215 <sup>a</sup>	2	.016
Likelihood Ratio	8.491	2	.014
Linear-by-Linear Association	7.911	1	.005
N of Valid Cases	106		

a. 0 cells (.0%) have expected count less than 5. The minimum expected count is 11.53.

Symmetric	Measures

		Value	Asymp. Std. Error <sup>a</sup>	Approx. T <sup>b</sup>	Approx. Sig.
Ordinal by Ordinal	Ordinal by Ordinal Kendall's tau-b		.086	-3.022	.003
	Spearman Correlation	276	.091	-2.924	.004°
Interval by Interval	Pearson's R	274	.091	-2.911	.004°
N of Valid Cases		106			
a. Not assuming	the null hypothesis.				
b. Using the asyn	nptotic standard error ass	uming the n	ull hypothesis.		
c. Based on norm	al approximation.				

a. Not assuming the null hypothesis.

b. Using the asymptotic standard error assuming the null hypothesis.

c. Based on normal approximation

#### Tests of Independence of Homework and Energy confounded by Gender

There is significant evidence that homework level and energy level are dependent (associated) in both the partial table for males and the partial table for females .....so we say:

#### Homework and Energy and conditionally dependent, given gender.

There is a conditional association of homework and energy for males.

#### There is a conditional association of homework and energy for females.

There is significant evidence that homework and energy are dependent when we look at the marginal table (that ignores male and female and puts all the counts together).

We have marginal dependence (association)of homework and energy.

#### Test of Linear Association of Homework and Energy confounded by Gender

There is significant evidence of a difference in median for the low and high homework groups when we condition on females. The negative sign on r indicates that low homework performers have a higher median energy rating than high homework performers for the female group.

There is significant evidence of a difference in medians for the low and high homework groups when we condition on males. The negative sign on r indicates that low homework performers have a higher median energy rating than high homework performers for the male group.

There is significant evidence of a difference in medians for the low and high homework groups when we look at the entire group of males and females together (the marginal result). The negative sign on r indicates that low homework performers have a higher median energy rating than high homework performers for the marginal situation.

age J

#### **Generalized Linear Models (GLM)**

- 1. Relate Explanatory and Response Variables
- 2. Models describe pattern of Association
- 3. Parameters describe nature/strength of Association
- 4. Test for Association: based on sample data while controlling for confounding variables
- 5. Can use models for Prediction

GLMs have both a random and a systemic component.

<u>Random Component:</u>This identifies that response variable Y and its distribution. The n observations Y1, Y2, ...Yn are independent. The Ys have categorical outcomes. <u>Binary data</u>: 2 possible outcomes (Success or Failure – coded 1 and 0 in raw data) Example: Y = 1 or Y = 0 (Y is Bernouilli or Binomial with n =1) <u>Count data</u>: Counts are outcomes Example: Y = number of successes in a given time interval or a specified region of space (Y is Poisson)

<u>Systematic Component</u>: This specifies how explanatory variables are entered in linear way into model with a "linear predictor"  $\alpha + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_k x_k$ 

<u>Link Functions</u>: We will consider various link functions that link random and systematic components using a function  $g(\mu)$ , where  $\mu$  is the mean of Y.

Model	$E(\mathbf{Y}) = \mu$ is expressed in terms of observed x values and unknown parameters	Link Function
Linear EXAMPLE: Binary Y, X=x	$\mu = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$ $\mu = \pi(x) = \alpha + \beta x$	Identity: $g(\mu) = \mu$
Logistic Regression <i>EXAMPLE:</i> <i>Binary Y, X=x</i>	$Log(\mu/(1-\mu)) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$ $\mu = \pi(x) = \frac{exp(a + \beta x)}{1 + exp(a + \beta x)}$	Logit: $g(\mu) = \log(\mu/(1-\mu))$
Probit Regression <b>EXAMPLE:</b> <i>Binary Y, X=x</i>	Probit( $\mu$ ) = F <sup>-1</sup> ( $\mu$ ) = $\alpha$ + $\beta_1 x_1$ + $\beta_2 x_2$ ++ $\beta_k x_k$ where $\mu$ = P(Z <= $\alpha$ + $\beta_1 x_1$ + $\beta_2 x_2$ + + $\beta_k x_k$ ) = F ( $\alpha$ + $\beta_1 x_1$ + $\beta_2 x_2$ + + $\beta_k x_k$ ) for the standard normal cdf $\mu$ = $\pi(x)$ = P(Z<= $\alpha$ + $\beta x$ )	Probit: $g(\mu)$ = inverse cumulative distribution of the standard normal distribution
Loglinear EXAMPLE: <i>Poisson Y, X =x</i>	$Log(\mu) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$ $\mu(x) = exp(\alpha + \beta x) = exp(\alpha)exp(\beta x)$	Log: $g(\mu) = \log(\mu)$

Exponential Family of Probability Distributions

These distributions can be written in a certain form, and include the normal, binomial, and poisson distribution.

Page 21

Each of these distributions has a "natural" parameter upon which its probabilities depend Normal:  $\mu$ ,  $\mu$  = mean Binomial: ln( $\pi$  /1-  $\pi$ ),  $\pi$ = success probability Poisson: ln  $\mu$ ,  $\mu$  = mean

The link function for "natural" parameter is called the "canonical" link. We use Maximum Likelihood Estimators with GLMs.

#### Example: HSB (High School and Beyond): N = 600 US high school seniors

For the next example, we will use a sample of 600 high school seniors, taken from the Large scale Longitudinal Study conducted by National Opinion Research Center (1980) under contract with National Center for Education Statistics. It includes the following variables. ID Student Identification Number (3 digits) SEX 1-Male 2-Female RACE (Race or Ethnicity) 1-Hispanic 2-Asian 3 - Black 4- White SES Socioeconomic Status 1-Low 2-Medium 3-High SCTYP School Type 1-Public 2-Private HSP High School Program 1-General 2-Academic Preparatory 3- Vocational/technical LOCUS Locus of Control Standardized to mean 0 and st dev 1 CONCPT Self Concept Standardized to mean 0 and st dev 1 MOT Motivation Average of 3 motivation items CAR Career Choice 1-Clerical 2-Craftsperson 3- Farmer 4-Homemaker 5-Laborer 6-Manager 7-Military 8-Operative 9-Professional<sup>1</sup>, 10-Professional<sup>2</sup>, 11-Proprietor 12-Protective 13-Sales 14-School 15-Service 16-Technical 17-Not working RDG Reading T-score standardized to mean 50 and SD 10 WRTG Writing T-score standardized to mean 50 and SD 10 MATH Math T-score standardized to mean 50 and SD 10 SCI Science T-score standardized to mean 50 and SD 10 CIV Civics T-score standardized to mean 50 and SD 10 1-Accountatn, Nurse, Engineer, Librarian, Writer, Social Worker, Athlete, Politician, etc 2 - Clergy, Dentist, Physician, Lawyer, Scientist, College Teacher, etc

#### High School and Beyond (HSB Data)

Below we present SPSS commands to calculate the linear equation for our GLM of interest, and discuss and interpret results. In SPSS terminology, we "reference" the "non-academic" schools, as our level of most interest is "academic" schools, and we wish  $\pi(x)$ , the probability of success to be the probability of an academic school. We will look at two possible ways to input SPSS commands to obtain information of interest. Note that it is important to choose the correct reference category for our chosen dependent response variable.

#### **Response:** BINARY

Y: Non-Academic High School – 0 Academic High School – 1\* OR Program: Academic High School – 0 Non-Academic High School -1) (\*It is intuitive to think of higher number 1 being assigned to level of most interest where our  $\pi(x)$  will be the probability of our "success". We will look at how SPSS handles this situation where the number 0 has been assigned to the level of most interest.)

**Explanatory:** (x = total of scores on 5 standardized achievement tests (Reading, Writing, Math, Science, Civics). This ranges from 164.7 to 350.

#### Random Component for Response Variable Y=program: Binomial Systemic Component for Explanatory Variable X=x: Linear Link Function: Identity $g(\mu) = \mu = (\alpha + \beta x)$

Analysis>Generalized Linear Models>Generalized	Analysis>Generalized Linear Models>Generalized
Linear Models	Linear Models
Type of Model Tab:	Type of Model Tab:
Custom:	Custom:
Distribution: Binomial	Distribution: Binomial
Link Function: Identity	Link Function: Identity
Response Tab:	Response Tab:
Dependent Variable: y(change Binary reference	Dependent Variable: program (default Binary
category to first – lowest value)	reference category is last – highest value)
	(Academic High School $-0$ or Non-Academic High
(Academic High School – 1 or Non-Academic	School -1)
High School -0)	Predictor Tab:
Predictor Tab:	x in the covariate box
x in the covariate box	Model Tab:
Model Tab:	x as a main effect
x as a main effect	✓ intercept in model
✓ intercept in model	All other tabs:
All other tabs:	left as defaults.
left as defaults.	

Page 22

#### Warnings

The maximum number of step-halvings was reached but the log-likelihood value cannot be further improved. Output for the last iteration is displayed.

The GENLIN procedure continues despite the above warning(s). Subsequent results shown are based on the last iteration. Validity of the model fit is uncertain.

There is at least one invalid case in the last iteration. A case is invalid if there are errors in computing the inverse identity link function, the log-likelihood, the gradient, or the Hessian matrix in the iterative process. Only the iteration history is displayed.

Execution of this command stops.

#### Iteration History

Iteration	Update Type	Number of Step-halvings	Log Likelihoodª	(Intercept)	x	(Scale)	Number of Invalid Cases <sup>b</sup>
0	Initial	0	-327.270	968450	.005700	1	4
1	Scoring	0	-326.932	848214	.005253	1	
2	Newton	0	-326.789	897262	.005432	1	2
3	Newton	1	-326.788	893459	.005419	1	2
4	Newton	5	-326.788°	893459	.005419	1	2

Redundant parameters are not displayed. Their values are always zero in all iterations. Dependent Variable: program

Model: (Intercept), x

a. The full log likelihood function is displayed.

b. These cases are invalid because there are errors in computing the inverse identity link function, the loglikelihood, the gradient, or the Hessian matrix in the iterative process. Invalid cases are excluded from iterations in which they cause computational errors.

c. The maximum number of step-halvings was reached but the log likelihood value cannot be further improved.

Intercept: Proportion of children in an academic high school with a total score of x=0 is -.893 (not sensible) Slope: For every increase of 1 point in x, the probability of being in an academic high school increased by .0054

Discussion and Interpretation:

Problems with using the identity link function.

1. The model takes values over the entire real line, but probabilities are always between 0 and 1. But, with small

or large enough x,  $\alpha + \beta x$  will either be smaller than 0 or greater than 1.

2. ML estimators of  $\alpha$  and  $\beta$  must be found by iterative methods as no formulas exist

3. The Fitting process fails when  $\hat{\pi}(x) = \hat{\alpha} + \hat{\beta}x$  outside of 0 to 1 range

The fitted line we obtained was:  $\hat{\pi}(x) = -0.893 + 0.0054x$ 

Note that these commands yield a column of  $\hat{\pi}(x)$ , and assign  $\hat{\pi}(x) = 0$  for observation 133 (where x = 164.7 is smallest x)  $\hat{\pi}(x) = 1$  for observation 84, (where x = 350 is largest x).

#### Logistic Regression Model

#### SPSS: Random Component for Response Variable Y=program: Binomial Systemic Component for Explanatory Variable X=x: Linear

Link Function: Logit:  $g(\mu) = ln(\frac{\mu}{1-\mu}) = (\alpha + \beta x)$ 

Recall, for the Logit Link function,

 $Log(\frac{\pi(x)}{1-\pi(x)}) = \alpha + \beta x \text{ OR } \pi(x) = \frac{\exp(a+\beta x)}{1+\exp(a+\beta x)}$ 

An advantage here is that  $\pi(x) = \frac{\exp(a + \beta x)}{1 + \exp(a + \beta x)}$  falls between 0 and 1 for all values of x.

#### SPSS COMMANDS FOR GLM APPROACH:

Analysis>Generalized Linear	Analysis>Generalized Linear Models>Generalized Linear				
Models>Generalized Linear Models	Models				
Type of Model Tab:	Type of Model Tab:				
Custom:	Custom:				
Distribution: Binomial	Distribution: Binomial				
Link Function: Logit	Link Function: Logit				
Response Tab:	Response Tab:				
Dependent Variable: y	Dependent Variable: program (default Binary reference				
(change Binary reference category to first –	category is last – highest value)				
lowest value)	(Academic High School $-0$ or Non-Academic High School -				
(Academic High School – 1 or Non-Academic	1)				
High School -0)	Predictor Tab:				
Predictor Tab:	x in the covariate box				
x in the covariate box	Model Tab:				
Model Tab:	x as a main effect				
x as a main effect	✓ intercept in model				
✓ intercept in model	All other tabs:				
All other tabs:	left as defaults.				
left as defaults.					

#### ONLY OUTPUT FOR THE DEPENDENT VARIABLE PROGRAM IS INCLUDED

Mod	lel Information	Goo	dness of Fit <sup>b</sup>		
			Value	df	Value/df
ependent variable	program	Deviance	612.542	520	1.178
robability Distribution	Binomial	Scaled Deviance	612.542	520	
nk Eurotion	Logit	Pearson Chi-Square	530.969	520	1.021
	Logi	Goodness of Fit <sup>®</sup> Value     df     Value/df       Deviance     612.542     520     1.178       Scaled Deviance     612.542     520     1.021       Pearson Chi-Square     530.969     520     1.021       Scaled Pearson Chi-Square     530.969     520     1.021       Scaled Pearson Chi-Square     530.969     520     1.021       Log Likelihood <sup>a</sup> -325.509     -325.509			
a. The procedure mo treating 1.00 as the r	dels .00 as the response, eference category.	Log Likelihood <sup>a</sup>	-325.509		
		Akaike's Information Criterion (AIC)	655.019		
		Finite Sample Corrected AIC (AICC)	655.039		
		Bayesian Information Criterion (BIC)	663.813		
		Consistent AIC (CAIC)	665.813		
		Dependent Variable: progra Model: (Intercept), x	am		
		a. The full log likelihood computing information c	function is dis riteria.	played and	used in
		b. Information criteria an	e in small-is-b	etter form.	

Page∠

Case Pro	ocessing S	umma	згу												
	N	Pe	rcent	1		Continuous Variable Information									
Included	600	10	0.0%	1				N	Mini	ուտ	Maximu	ın Me	an	Std. De	viation
Excluded	0		.0%			( ovariate	χ	កព	) 1,	64 70	960	00 269	447	4144844	
Total	600	10	0.0%			Tests of Model Effects									
				1					٦	Fype I	II				
	Omnibus 1	'estª				Source		Wald Cl Squar	hi- e	(	lf	Sig.			
Likelihood Ratio Chi-					(Interce	ept)	103	.104		1	.00	5			
Square	C	ſ	Sig	l.		x		106.411			1	.00			
139.0	139.082 1 .000				Depen	dent V	ariable: pr	ogram							
Dependent Model: (Inte	Variable: p rcept), x	rograr	n			Model:	(Interc	ept), x	Para	ameter l	Estimates				
a. Comp	ares the fit	ted ma	odel						95% Wal	d Confid	ence Interva	al	Нуро	thesis Test	
againsti	the interce	ot-only	model			Parameter	в	Std. Error	Lowe	r	Upper	Wald C Squa	hi- e	df	Sig.
	Categorica	l Variab	le Inforr	nation		(Intercept)	-7.05	5.6948	-8.	.417	-5.69	3 10:	.104	1	.000
				Ν	Percent	x (Scale)	.027 14	7 .0027 a		.022	.03	3 10	i.411	1	.000
Dependent Va	riable pro	gram	.00	308	51.3%	Dependent	l /ariable: p	program							
			1.00	292	48.7%	a. Fixed a	cept), x t the displ	laved value.							
	β detern	nines	the rate of	of incr	rease	or de	crease o	f the	e S cur	ve of					
						the logis	stic re	gression	mode	el.					
							Logit $(\widehat{\pi(x)}) = -7.055 + 0.027x$ is estimated equation								
						Interpretation: Since $0.027$ is > 0, the estimated probability of									
						attending an academic school increases as x (total score)									
						increase	s.								

#### SPSS COMMANDS FOR REGRESSION: BINOMIAL APPROACH

#### Random Component for Response Variable Y=program: Binomial

Systemic Component for Explanatory Variable X=x: Linear

Link Function: Logit:  $g(\mu) = ln(\frac{\mu}{1-\mu}) = (\alpha + \beta x)$ 

Analysis>Regression>Binary Logistic

Dependent: y

Covariate(s):x

Leave other defaults

Ok

<u>This automatically References</u> : (your lowest coded category) For y: "non-academic schools"

#### OUTPUT

#### **Omnibus Tests of Model Coefficients**

		Chi-square	df	Sig.
Step 1	Step	139.082	1	.000
	Block	139.082	1	.000
	Model	139.082	1	.000

#### **Model Summary**

Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1	692.268 <sup>a</sup>	.207	.276
<b>D</b>		• • • •	1

a. Estimation terminated at iteration number 4

because parameter estimates changed by less than .001.

#### Variables in the Equation

		В	S.E.	Wald	df	Sig.	Exp(B)
Step 1a	Х	.027	.003	106.411	1	.000	1.028
	Constant	-7.055	.695	103.104	1	.000	.001

a. Variable(s) entered on step 1: x.

We can readily find the coefficients for the equation Logit  $(\widehat{\pi(x)}) = -7.055 + 0.027x$  above.

#### **Probit Regression Model**

Random Component for Response Variable Y=program: Binomial Systemic Component for Explanatory Variable X=x: Linear Link Function: Probit  $(F(\alpha + \beta x)) =$  Probit  $(P(Z \le (\alpha + \beta x)))$ Probit $(\pi(x)) = \alpha + \beta x$  where  $\pi(x) =$  left tail probability for standard normal distribution with z score  $\alpha + \beta x$  OR  $\pi(x) = P(Z \le \alpha + \beta x) = F(\alpha + \beta x)$  on  $-\infty \le x \le \infty$  for standard normal distribution Transforms  $\pi(x)$  so that the regression curve for  $\pi(x)$  has the appearance of the normal cdf Advantage:  $\pi(x)$  falls between 0 and 1 for all x.

Page Z

Analysis>Generalized Linear Models>Generalized	Analysis>Generalized Linear Models>Generalized
Linear Models	Linear Models
Type of Model Tab:	Type of Model Tab:
Custom:	Custom:
Distribution: Binomial	Distribution: Binomial
Link Function: Probit	Link Function: Probit
Response Tab:	Response Tab:
Dependent Variable: y	Dependent Variable: program (default Binary
(change Binary reference category to first - lowest	reference category is last – highest value)
value)	(Academic High School $-0$ or Non-Academic High
(Academic High School – 1 or Non-Academic High	School -1)
School -0)	Predictor Tab:
Predictor Tab:	x in the covariate box
x in the covariate box	Model Tab:
Model Tab:	x as a main effect
x as a main effect	✓ intercept in model
✓ intercept in model	All other tabs:
All other tabs:	left as defaults.
left as defaults.	

#### OUTPUT is for dependent variable program only

Mod	el Information
Dependent Variable	program <sup>a</sup>
Probability Distribution	Binomial
Link Function	Probit

a. The procedure models .00 as the response, treating 1.00 as the reference category.

#### **Case Processing Summary**

	N	Percent
Included	600	100.0%
Excluded	0	.0%
Total	600	100.0%

#### Omnibus Test<sup>a</sup>

Likelihood Ratio Chi- Square	df	Sig.
138.947	1	.000

Dependent Variable: program Model: (Intercept), x

a. Compares the fitted model against the intercept-only model.

Goodness of Fit <sup>b</sup>						
Value df Va						
Deviance	612.677	520	1.178			
Scaled Deviance	612.677	520				
Pearson Chi-Square	532.430	520	1.024			
Scaled Pearson Chi- Square	532.430	520				
Log Likelihood <sup>a</sup>	-325.577					
Akaike's Information Criterion (AIC)	655.153					
Finite Sample Corrected AIC (AICC)	655.174					
Bayesian Information Criterion (BIC)	663.947					
Consistent AIC (CAIC)	665.947					

Dependent Variable: program Model: (Intercept), x

a. The full log likelihood function is displayed and used in computing information criteria.

b. Information criteria are in small-is-better form.

#### **Continuous Variable Information**

	N	Minimum	Maximum	Mean	Std. Deviation
Covariate x	600	164.70	350.00	259.9447	40.44849

Categorical Variable Information					
			N	Percent	
Dependent Variable	program	.00	308	51.3%	
		1.00	292	48.7%	
Total 600 100.0%					

Parameter Estimates							
			95% Wald Con	35% Wald Confidence Interval Hypothesis Te			
Parameter	в	Std. Error	Lower	Upper	Wald Chi- Square	df	Sig.
(Intercept)	-4.252	.3932	-5.023	-3.482	116.965	1	.000
х	.017	.0015	.014	.019	121.093	1	.000
(Scale)	1ª						
Dependent Variable: program Model: (Intercept), x							
a Fixed a	t the display	ed value					

Tests of Model Effects

	Type III					
Source	Wald Chi- Square	df	Sig.			
(Intercept)	116.965	1	.000			
х	121.093	1	.000			
Dependent Variable: program						

Model: (Intercept), x

 $\beta$  determines the rate of increase or decrease of the S curve of the probit regression model.

Probit( $\widehat{\pi(x)}$ ) = -4.252 + 0.017x is estimated equation (see highlighted parameter estimate output on following output slide)

Interpretation: Since 0.017 is > 0, the higher the total score, the higher a probability is that a person attends an academic high school program.

#### Example 29: Aids/Australia

Model: Poisson Loglinear Model



Random Component for Response Variable Y = # deaths: Poisson Systemic Component for Explanatory Variable X = x =time: Linear Link Function:  $Ln(\mu) = \alpha + \beta x$ 

#### COMMANDS

Analysis>Generalized Linear	To put the equation on the scatterplot,		
Models>Generalized Linear Models	Double click graph so pops up		
Type of Model Tab:	Right click on graph		
– Custom:	Select: Add Reference Line from Equation		
<ul> <li>Distribution: Poisson</li> </ul>	Custom Equation: Type $exp(0.340 + 0.257*x)$		
<ul> <li>Link Function: Log</li> </ul>	Apply		
Response Tab:			
<ul> <li>Dependent Variable: Y</li> </ul>	$\hat{\mu}(x) = \exp(0.340 + 0.257x)$ is estimated equation and		
Predictor Tab:	can be read from the output below		
<ul> <li>x in the covariate box</li> </ul>			
Model Tab:			
<ul> <li>x as a main effect</li> </ul>			
<ul> <li>include intercept in model</li> </ul>			
All other tabs: left as defaults.			

#### <u>OUTPUT</u>

Model Information				Good	ness of Fit <sup>b</sup>	
Dependent Variable	Y		ſ		Value	
Probability Distribution	Poisson		ſ	Deviance	29.654	-
Link Function	Log			Scaled Deviance	29.654	
				Deersen Ohi Onvers	20.047	

#### **Case Processing Summary**

	N	Percent
Included	14	100.0%
Excluded	0	.0%
Total	14	100.0%

#### Omnibus Test<sup>a</sup>

Likelihood Ratio Chi- Square	df	Sig.
177.619	1	.000

Dependent Variable: Y Model: (Intercept), x

a. Compares the fitted model against the intercept-only model.

#### Continuous Variable Information

		N	Minimum	Maximum	Mean	Std. Deviation
Dependent Variable	Y	14	.00	45.00	15.6429	14.98516
Covariate	х	14	1.00	14.00	7.5000	4.18330

Goodness of Fit <sup>b</sup>							
	Value	df	Value/df				
Deviance	29.654	12	2.471				
Scaled Deviance	29.654	12					
Pearson Chi-Square	28.847	12	2.404				
Scaled Pearson Chi- Square	28.847	12					
Log Likelihood <sup>a</sup>	-41.290						
Akaike's Information Criterion (AIC)	86.581						
Finite Sample Corrected AIC (AICC)	87.672						
Bayesian Information Criterion (BIC)	87.859						
Consistent AIC (CAIC)	89.859						

Dependent Variable: Y Model: (Intercept), x

a. The full log likelihood function is displayed and used in computing information criteria.

b. Information criteria are in small-is-better form.

#### **Tests of Model Effects**

	Type III						
Source	Wald Chi- Square	df	Sig.				
(Intercept)	1.828	1	.176				
х	135.477	1	.000				

Dependent Variable: Y Model: (Intercept), x

Parameter Estimates										
			95% Wald Cont	idence Interval	Нуро	thesis Test			95% Wald Conf for Ex	idence interval p(B)
Parameter	в	Std. Error	Lower	Upper	Wald Chi- Square	df	Sig.	Exp(B)	Lower	Upper
(Intercept)	.340	.2512	153	.832	1.828	1	.176	1.404	.858	2.298
х	.257	.0220	.213	.300	135.477	1	.000	1.292	1.238	1.349
(Scale)	1ª									
Dependent Variable: Y Model: (Intercept), x										
a. Fixed at	a. Fixed at the displayed value.									

 $_{\rm Page}29$ 

#### LOGLINEAR MODEL

Log( $\mu$ ) =  $\alpha$  +  $\beta$ x OR  $\mu$  = exp( $\alpha$  +  $\beta$ x) = e<sup> $\alpha$ </sup>e<sup> $\beta$ x</sup> A one unit increase in x has a multiplicative impact of e<sup> $\beta$ </sup> on  $\mu$ . If  $\beta$  >0, then e<sup> $\beta$ </sup> >1, and  $\mu$  increases as x increases. If  $\beta$  < 0, then  $\mu$  decreases as x increases.  $\mu(x) = e^{0.340 + 0.257x}$  is estimated equation (from output) Interpretation: The number of AIDS deaths in OZ over a 3 month period was in average e<sup>0.257</sup> = 1.29 times higher than in the previous 3 month period.

Inference re Model Parameter  $\beta$ 

Ho:  $\beta = 0$  versus Ha:  $\beta \neq 0$ 

Assumptions: Random samples large enough that ML estimators approximately normal

	Possible Test Statistics	P-value				
Wald Z	$z_o = \hat{\beta}/SE(\hat{\beta})$ (is Z under Ho)	$2P(Z >  z_o )$				
Wald $\chi^2$	$\chi_0^2 = z_0^2 = (\hat{\beta}/SE(\hat{\beta}))^2$ , df=1 (is $\chi^2$ under Ho)	$P(\chi^2 > \chi^2_o)$				
Likelihood ratio	$-2\ln(l_0/l_1) = -2(L_0-L_1), df=1$ (is $\chi^2$ under Ho)	$P(\chi^2 > \chi^2_o)$				
$\hat{\beta} =$ maximum likelihood	$\hat{\beta}$ = maximum likelihood estimate of $\beta$					
$l_0$ = maximized value of likelihood function under Ho						
$l_1$ = maximized value of likelihood function when $\beta$ need not equal 0						
$L_0 = \ln(l_1)$ and $L_1 - \ln(l_0)$	are the maximized log-likelihood functions					

From output:

Wald  $\chi^2$   $\chi_o^2 = 135.477$ , df=1  $P(\chi^2 > \chi^2_o) < 0.001$ 

Reject Ho. Have sufficient evidence that time has an effect on the number of AIDS deaths in Oz.

#### Generalization: Compare Model of Interest M to Saturated Model

Ho:proposed model M fits as well as saturated model S vs Ha:Ho not correct Ho:all parameters in model S that are not in model M = 0 (said another way) Assumptions: Random samples, large samples

Rejection of Ho suggests that we need a full saturated model (model of "perfect fit" with one parameter per "measurement"). Being able to not reject Ho suggests that model M might suffice. Rejection of Ho means model M has significantly worse fit than model S. Therefore we hope test is **not** significant!

We test this with a:

GOODNESS OF FIT TEST BASED ON **DEVIANCE =** -2(L<sub>M</sub>-L<sub>S</sub>)

Likelihood ratio test (Deviance test)	Deviance <sub>0</sub> = -2(L <sub>M</sub> -L <sub>S</sub> )	For some GLMs*, test statistic is $\chi^2$ with df=difference in number of parameters in models S and M when Ho is true *holds for Poisson loglinear, for example	$P(\chi^2 > \chi^2_o)$ when Ho is true			
$l_M = maximized$	$l_M$ = maximized value of likelihood function under proposed model M					
$l_{S =} maximized$	$l_S$ = maximized value of likelihood function under saturated model S					
$L_M = maximized$	$L_M$ = maximized log likelihood of model of interest and $L_S$ = maximized log likelihood of					
saturated model	saturated model					

Measures of Deviance/df "close" to 1 indicate good model fit.

Aids in Australia: Ho: model of interest fits as well as (full) saturated model (Poisson loglinear)

Goodness of Fit test	$-2(L_M-L_S) = 29.654, df = 14-2 = 12$	Deviance/df = 29.654/12 = 2.471
on output		

We would reject Ho. The saturated model fits significantly better than the proposed model.

#### Generalization: Compare two models, M<sub>0</sub> is nested in M<sub>1</sub>

Ho:model  $M_0$  fits as well as model  $M_1$  vs Ha:Ho not correct Ho:all parameters in model  $M_1$  but not in model  $M_0=0$  (said another way)

Assumptions: Random samples, large samples

Likelihood ratio	$-2(L_0-L_1) = Deviance_0 - Deviance_1 \text{ is } \chi^2 \text{ with df}=difference in number of parameters in models M_0 and M_1$	$P(\chi^2 > \chi^2_o)$ when Ho is true
$L_0 = maximized$ Deviance <sub>0</sub> – De	l log likelihood of model Mo and $L_1$ = maximized log likelihood of Model M viance <sub>1</sub> = -2(L <sub>0</sub> -L <sub>S</sub> ) - 2(L <sub>1</sub> -L <sub>S</sub> ) = -2(L <sub>0</sub> -L <sub>1</sub> )	M1

Rejection of Ho with a "relatively" large test statistic means and small p-value means that model  $M_0$  fits poorly in comparison to model  $M_1$ , and suggests the larger model is significantly more appropriate.

Aids in Australia: Ho:Intercept only model fits as well as model including time

Likelihood ratio	$-2\ln(l_0/l_1) = -2(L_0-L_1) = 177.619, df = 2-1=1$	$P(\chi^2 > \chi^2_{o}) < 0.001$
	(See Omnibus test on Output)	

There is sufficient evidence that the model including time fits better than the intercept only model.

#### **Residuals**

These are "standardized" differences between the observed and estimated values, and should be investigated when they exceed 2 in absolute value.

Pearson's Residuals  $ei = \frac{yi - \hat{\mu}i}{\sqrt{Var}(yi)}$ 

If you group the x variable (so it is categorical), these are approximately normal.

If you don't group the x variable, these are not approximately normal.

The Standardized Pearson's Residuals  $ei = \frac{yi - \hat{\mu}i}{SE}$  are approximately normal, where SE is estimated standard error of residual  $yi - \hat{\mu}_i$ .

Deviance Residuals measure how much individual measurements add to the deviance of a model.

Random Component for Response Variable Y=: Poisson Systemic Component for Explanatory Variable X= Time: Linear Link Function: Log

The following commands can be used to obtain the necessary residuals.

#### **COMMANDS**

Analysis>Generalized Linear Models>Generalized Linear Models Type of Model Tab: Counts: Poisson loglinear Response Tab: Dependent Variable: Y Predictor Tab: x in the covariate box Model Tab: x as a main effect include intercept in model Save Tab: ✓Pearson residual ✓ Standardized Pearson residual ✓ Standardized deviance residual All other tabs:

left as defaults.

After save: the datasheet columns are as follows

Х	Y	MeanPredicted	PearsonResidual	DevianceResidua	l StdPearsonResidual
1.00	.0	1.82	-1.347	-1.905	-1.417
2.00	1.00	2.35	879	993	928
3.00	2.00	3.03	593	632	627
4.00	3.00	3.92	464	484	492
5.00	1.00	5.06	-1.806	-2.210	-1.916
6.00	4.00	6.55	995	-1.073	-1.054
7.00	9.00	8.46	.186	.184	.196
8.00	18.00	10.93	2.137	1.953	2.249
9.00	23.00	14.13	2.359	2.161	2.475
10.00	31.00	18.26	2.980	2.708	3.128
11.00	20.00	23.60	742	762	785
12.00	25.00	30.51	997	-1.029	-1.084
13.00	37.00	39.43	386	391	449

NOTE LARGER RESIDUALS AT x = 8, 9, AND 10. At that time, it appears that other factors were involved.

The HSB data will be used to create several GLMs for investigation. Commands, output, and interpretations will be offered.

The following commands allow for the creation of a logit expression when Y is the response variable and the explanatory variable is X, the academic score.

Example 35: (Using Regression path)	Respon	Response: Y – (0-non-academic, 1-academic)						
NOTE: DEFAULT SUCCESS	Explanatory:							
PROBABILITY IS THE ONE WITH	X – academic score (total of reading, writing, math, science,							
THE HIGHEST LEVEL IN THE	civics)							
RESPONSE VARIABLE		Model Summary						
(ACADEMIC)		-21	ວດ	Cox	& Snell R	Nage	lkerke R	
ANALYZE>REGRESSION>BINARY	Star	1:11:1-	55 1	COA		ruge		
LOGISTIC	Step	likelin	000	2	quare	2	quare	
Logistic Regression Screen	1	69	92.268 <sup>a</sup>		.207		.276	
-Move y to Dependent box	a. Estir	nation te	rminate	ed at i	iteration n	umber	: 4	
-Move x to Covariate box	because	e parame	ter esti	mates	s changed	by les	s than	
Covariate box will read	.001	1			U			
X								
-Method: Enter								
(Use Save and Options buttons as	Variables in the Equation							
needed)			В	S.E.	Wald d	f Sig.	Exp(B)	
- UK								
	Step 1 <sup>a</sup>	Х	.027	.003	106.411	1 .000	1.028	
	Step 1 <sup>a</sup>	x Constant	.027 -7.055	.003 .695	106.411 103.104	1 .000 1 .000	1.028 .001	
	Step 1 <sup>a</sup> a. Varia	x Constant able(s) er	.027 -7.055 ntered o	.003 .695 on ste	106.411 103.104 p 1: x.	1 .000 1 .000	1.028 .001	
	Step 1 <sup>a</sup> a. Varia	x Constant able(s) er	.027 -7.055 ntered (	.003 .695 on ste	106.411 103.104 p 1: x.	1 .000 1 .000	1.028 .001	
	Step 1 <sup>a</sup> a. Varia $\pi(\mathbf{x})$ , su	x Constant able(s) en ccess pr	.027 -7.055 ntered o <b>ob is p</b>	.003 .695 on ste <b>rob o</b>	106.411 103.104 p 1: x. f attendi	1 .000 1 .000	1.028 .001	chool
	Step 1 <sup>a</sup> a. Varia $\pi(x)$ , su Equatio	x Constant able(s) er ccess pr n: logit(i	.027 $-7.055$ intered of the second secon	.003 .695 on ste <b>rob o</b> -7.05	106.411 103.104 pp 1: x. of attendi 55 + 0.027	1 .000 1 .000 ng aca	1.028 .001	chool
	Step 1 <sup>a</sup> a. Varia $\pi(\mathbf{x})$ , su Equatio Slope: a	x Constant able(s) en ccess pr n: logit(2 a 1 point	.027 -7.055 intered of <b>ob is p</b> $\widehat{\boldsymbol{\tau}(\boldsymbol{x})} =$ increas	.003 .695 on ste <b>rob o</b> -7.05 se in t	106.411 103.104 p 1: x. <b>f attendi</b> 55 + 0.027 he total te	1 .000 1 .000 ng aca 7x est scor	1.028 .001	<b>chool</b> iplies the
	Step 1 <sup>a</sup> a. Varia $\pi(\mathbf{x})$ , su Equation Slope: <i>a</i> odds of	x Constant able(s) er ccess pr n: logit(a a 1 point a studen	.027 -7.055 ntered o <b>ob is p</b> $\overline{\tau(x)} =$ increas t being	.003 .695 on ste <b>rob o</b> -7.05 se in t in an	106.411 103.104 p 1: x. of attendi 55 + 0.027 he total te academic	1 .000 1 .000 ng aca $7_X$ est score c progr	1.028 .001	<b>chool</b> iplies the factor of
	Step 1 <sup>a</sup> a. Varia $\pi(\mathbf{x})$ , su Equatio Slope: a odds of $e^{0.027}=$	x Constant able(s) en ccess pr n: logit( <i>i</i> a 1 point a studen 1.02. Th	.027 -7.055 ntered o <b>ob is p</b> $\overline{\tau(x)} =$ increas t being ne odds	.003 .695 on ste <b>rob o</b> -7.05 se in t in an s of be	106.411 $103.104$ ap 1: x. <b>of attendi</b> $55 + 0.027$ he total teal         academic         eing in an	$\begin{array}{c c} 1 & .000 \\ \hline 1 & .000 \\ \hline 1 & .000 \\ \hline ng aca \\ 7x \\ est score \\ c program \\ acade \\ \end{array}$	1.028 .001	<b>chool</b> iplies the factor of gram increase

The following commands allow information necessary for the calculation of the Wald Z, Wald  $\chi^2$ , and the Likelihood Ratio test statistics to be obtained, in addition to the logit expression, and the predicted  $\pi(x)$  values.

Example 35: Commands for Inference Results: Response: Y – (0-non-academic, 1-academic) Explanatory: X – academic score (total of reading, writing, math, science, civics) NOTE: DEFAULT SUCCESS PROBABILITY IS THE ONE WITH THE HIGHEST LEVEL IN THE **RESPONSE VARIABLE (ACADEMIC)** ANALYZE>REGRESSION>BINARY LOGISTIC Logistic Regression Screen -Move y to Dependent box -Move x to Covariate box *Covariate box will read* x -Method: Enter (Use Save and Options buttons as shown below) -Ok Option box: Check "CI for exp(B)" – will provide 95% CI for  $e^{\beta}$  in the Variables in the Equation box in the output file Save Box: Check "Probabilities" – will provide a column of predicted values for  $\pi(x)$  in the data file

To obtain predicted values for  $\pi(x)$  when your x of interest is not one of the x values in the data, add the x value of interest to the bottom of the x column, in the row below the last row of data.

Example: For x = 260, SPSS provides 0.5162 as predicted value for  $\pi(260)$  (in the data file).

Omnibus Tests of Model Coefficients

		Chi-square	df	Sig.
Step 1	Step	139.082	1	.000
	Block	139.082	1	.000
	Model	139.082	1	.000

#### **Model Summary**

	-2 Log	Cox & Snell R	Nagelkerke R		
Step	likelihood	Square	Square		
1	692.268ª	.207	.276		

a. Estimation terminated at iteration number 4 because

parameter estimates changed by less than .001.

#### Variables in the Equation

							95% C.I.for EXP(E		
	В	S.E.	Wald	df	Sig.	Exp(B)	Lower	Upper	
Step 1 <sup>a</sup> x	.027	.003	106.411	1	.000	1.028	1.022	1.033	
Constant	-7.055	.695	103.104	1	.000	.001			

a. Variable(s) entered on step 1: x.

#### Example 35 (commands for inference results) (continued)

We can use the above information to perform inference:

Tests of significance for the model:  $logit(\pi(x) = \alpha + \beta x)$ :

Ho:  $\beta = 0$  versus Ha:  $\beta \neq 0$  , choose  $\alpha = 5\%$ 

Assumptions: Random samples large enough that ML estimators approximately normal

 $\hat{\beta}$  = maximum likelihood estimate of  $\beta$ 

coloured information is from output above

	Test Statistic	P-value
Wald Z	$z_o = \hat{\beta}/SE(\hat{\beta})$ (is Z under Ho) = $\sqrt{106.411} = 10.32$	$\begin{array}{l} 2P( Z>  z_{o}  ) \\ <\!\!0.001 \ (from \ tables) \end{array}$
Wald $\chi^2$	$\chi_o^2 = z_o^2 = (\hat{\beta}/SE(\hat{\beta}))2$ , df=1 (is $\chi^2$ under Ho) =106.411, df =1	$\frac{P(\chi^2 > \chi^2_o)}{< 0.001}$
Likelihood ratio	$-2\ln(l_0/l_1) = -2(L_0-L_1), df=1 \text{ (is } \chi^2 \text{ under Ho)}$ = 139.082, df =1	$\frac{P(\chi^2 > \chi^2_o)}{< 0.001}$

Reject Ho

At a significance level of 5%, we have sufficient evidence that the odds of being in an academic program are related to the test results. (That is; we have sufficient evidence to support the model  $M_1$ .)

Wald CI for  $\beta$ : 0.027 +/- 1.96(0.003) = (.022,.033) (by hand using SPSS output information) Wald CI for  $e^{\beta} = (1.022, 1.033)$ 

**Example 35:** Additional inference output needed for confidence intervals for predicted values of  $\pi(x)$  can be found when one uses the <u>Generalized Linear Model</u> path to perform your binary logistic regression. Response: Y – (0-non-academic, 1-academic)

Explanatory:

X – academic score (total of reading, writing, math, science, civics)

Analysis>Generalized Linear Models>Generalized Linear Models

Type of Model Tab: Choose "Binary Logistic" button

Response Tab: Select y for the Dependent Variable *Binary button popup box: Change reference category to first (lowest value). This will ensure that your*  $\pi(x)$  *probability of success will be the probability of an academic school.* 

Predictors Tab: Select x for the Covariates box

Model Tab: Choose Type "Main Effects" and move the x variable into the Model box Estimation Tab: Leave as default

Statistics Tab: Leave all defaults AND also check "Include exponential parameter estimates" and "Covariance matrix for parameter estimate". (Also note that you can change the level of confidence of your CIs here, and that the defaults for the inference tests and the confidence interval types are Wald.) EM Means Tab: Leave as default

Save Tab: Check "Predicted value of mean of response, "Lower bound of confidence interval for mean of response", "Upper bound of confidence interval for mean of response", "Predicted Value of linear predictor", "Estimated standard error of predicted value of linear predictor".

Page

Click OK		Madel	Informati	0 <b>m</b>						
	• 1 1	Niodel	Informati	on		_				
Dependent Va	ariable	y. Din	omial							
Link Function	suibuu	Logit								
LIIK Function	1	L08	git		time 00 as the					
a. The proceed	are moc	iers 1.00	as the resp	onse, trea	ting .00 as the	e				
Telefence cate	gory.	Goo	dness of Fit	b						
Γ		000	Value	df	Value/df					
Deviance			612 542	520	1 1	78				
Scaled Deviand	ce		612.542	520	1.1	/0				
Pearson Chi-Sc	juare		530.969	520	1.0	21				
Scaled Pearson	Chi-Sq	uare	530.969	520						
Log Likelihood	1 <sup>a</sup>		-325.509							
Akaike's Inform	nation		655.019							
Criterion (AIC)	)	1 4 10	<55 0 <b>0</b> 0							
Finite Sample	Correcte	a AIC	655.039							
(AICC) Bayasian Infor	mation		663 812							
Criterion (RIC)	111at1011 )		003.013							
Consistent AIC	, C (CAIC)	)	665.813							
Dependent Va	riable	V	0001010							
Model: (Interc	cent) x	y								
a. The full log	likelih	ood fund	ction is disr	played and	used in					
computing inf	ormatic	on criteri	ia.	Juyea and						
b. Information	criteria	a are in s	small-is-be	tter form.						
	0	mnibus	Testa							
Likelihood										
Ratio Chi-										
Square		df		Sig.						
139.0	)82	1		U	.000					
Dependent Va	riable	V								
Model: (Interc	cent) x	y								
a. Compares f	he fitted	1 model	against the	intercent-	only					
model		a model	uguinist the	intercept	omy					
mouen	Tests	of Mod	el Effects							
			Type III							
	Wald	d Chi-								
Source	Sq	uare	df	Sig.						
(Intercept)	1	103.104	1	<u> </u>	.000					
X	1	106.411	1		.000					
Dependent Va	riable:	v								
Model: (Interc	cept), x	5								
	· · · · · ·			Para	meter Estimate	es				
			95% Wald	Confidence						
		~ •	Inte	erval	Нуро	thesis Te	st			
Doromotor	ъ	Std.	Louise	Unner	Wald Chi-	Af	<b>C</b> :~	$\mathbf{E}_{\mathbf{v}\mathbf{n}}(\mathbf{D})$		
(Intercept)	D	EITOF 6049	2 /17	opper 5.60	3 102 104	ui 1	51g.	Exp(B)		
(Intercept) X	.027	.0027	.022	.03	<b>3</b> 105.104	1	.000	1.028		
(Scale)	1 <sup>a</sup>									
							•			

 $_{\rm Page}36$ 

95% Wald Confidence Interval for Exp(B)

Lower

.000

1.022

Upper

.003 1.033
### **Profile Likelihood Confidence Intervals**

It appears that SPSS can calculate Profile Likelihood Confidence Intervals. You tell it on the Statistics tab to use a "Likelihood ratio" Chi-square statistics and a "Profile likelihood" Confidence Interval Type. It produces the output below. I'm not convinced it is really calculating the "Profile Likelihood" confidence intervals for  $\beta$ , given that they appear close to the "Wald" confidence intervals. However, the confidence intervals for the intercept are different, and if you look at further decimal places for the CI for  $\beta$ , they do differ. And when I ran some other models with more explanatory variables I did get some larger differences in CIs for various parameters for the two types of confidence intervals. **Output** 

Tests of Model Effects										
_	Type III									
	Likelihood	Likelihood								
	Ratio Chi-									
Source	Square	df	Sig.							
(Intercept)	133.838	1	.000							
X	139.082	1	.000							

Dependent Variable: y

Model: (Intercept), x

### **Parameter Estimates**

			95%	Profile					95% Profile	e Likelihood
			Likel	ihood					Confidence	Interval for
			Confiden	ce Interval	Hypothesis Test				Ex	p(B)
		Std.			Wald Chi-					
Parameter	В	Error	Lower	Upper	Square	df	Sig.	Exp(B)	Lower	Upper
(Intercept)	-7.055	.6948	-8.457	-5.730	103.104	1	.000	.001	.000	.003
х	.027	. <mark>0027</mark>	.022	.033	106.411	1	.000	1.028	1.023	1.033
(Scale)	1 <sup>a</sup>									

Dependent Variable: y

Model: (Intercept), x

a. Fixed at the displayed value.

### **Covariances of Parameter Estimates**

	(Intercept)	Х
(Intercept)	.48271	00183
Х	00183	7.04654E-6

**Profile Likelihood CI for β:** (.022,.033) **Profile Likelihood CI for**  $e^{\beta} =$  (1.022, 1.033)

Confidence Intervals for the Probabilities,  $\pi(x)$ 

From Data File (because of Save tab):

New Columns Appear (this is row 601, which has predicted value information for $\pi(260)$									
XBPredicte	KBPredicted SBStandardError MeanPredicted CIMeanPredictedLower CIMeanPredicte								
0.00	5	0.092	0.516	0.471	0.561				

95% CI for logit( $\pi(260)$ ) is 0.065 +/- 1.96(0.092) = (-0.115,0.244) (by hand, using computer information) Exponentiating, by hand, on the value of 0.065 gives .516 (for  $\pi(260)$ ) Exponentiating, by hand gives (0.471,0.561) as a 95% CI for  $\pi(260)$ )

 $\pi(260) = 0.516$  (from output) 95% CI for  $\pi(260)$  is (0.471, 0.561) (from output)

# SO WE DO HAVE A MATCH:

# ROUNDING PROBLEMS WHEN USING COMPUTER OUTPUT: BEWARE:

A suitable formula for calculating a CI for logit( $\pi(x)$ ) is:

logit( $\widehat{\pi(x)} \pm 1.96\sqrt{SE^2}$ , where SE<sup>2</sup> =  $Var(\widehat{\alpha}) + x^2 Var(\widehat{\beta}) + 2x Cov(\widehat{\alpha}, \widehat{\beta})$ 

If you calculate logit( $\widehat{\pi(260)}$ ) using the ROUNDED formula -7.055 + 0.027 provided in the SPSS output, you will obtain -0.035 as your logit( $\widehat{\pi(260)}$ ) If you calculate logit( $\widehat{\pi(260)}$ ) using 12 decimal places (which you can obtain by popping up the boxes from the SPSS output), you will obtain 0.065 as your logit( $\pi(260)$ ).

YOU WOULD NOT WANT TO USE THIS INCORRECT VALUE OF -0.35 IN CALCULATING YOUR CI BY HAND.

FURTHERMORE, IF YOU USE THE COVARIANCE OF PARAMETER ESTIMATES TABLE PROVIDED BY THE COMPUTER OUTPUT (SEE ABOVE) AND SUBSTITUTE THESE VALUES INTO THE EQUATION ABOVE, FURTHER ROUNDING ERROR WILL OCCUR AND YOUR INTERVAL WIDTH WILL BE INCORRECT.

HOWEVER, IF YOU CARRY 12 DECIMAL PLACES THROUGHOUT ALL YOUR BY HAND WORK, YOUR BY HAND ANSWERS WILL MATCH THE COMPUTER OUTPUT.

IT'S NICE WHEN COMPUTERS DO THINGS FOR US.

THIS PROBLEM OCCURS DUE TO THE SMALL MAGNITUDE OF THE logit( $\widehat{\pi(260)}$ )

Example 37:	Response: Y – (0-non-academic, 1-academic)	
NOTE: DEFAULT INDICATOR	Explanatory:	
REFERENCE CATEGORY FOR	SES (1-low, 2-middle, 3-high)	
EXPLANATORY VARIABLES IS	(express SES with 2 dummy variables)	
HIGHER NUMBERED LEVEL	$d_1 = 0 - not low, 1 - low$	
ANALYZE>REGRESSION>	$d_2 = 0 - not med, 1 - medium$	
BINARY LOGISTIC	Model Summary	
Logistic Regression Screen		
-Move y to Dependent box		

-Move SES to Covariate box <i>Covariate box will read</i> SES -Method: Enter -Click Categorical Button <u>Covariate Screen</u> -Move SES from Covariate box to categorical Covariates box <i>Categorical Covariate box will read</i> SES(indicator) Continue	Step 1 a. Est numb chang	-2 Log likelihood 779.39( imation ten er 4 becau ged by less	Co Sne d Squ D <sup>a</sup> rminate se para than .( 7 <b>ariabl</b>	x & ell R uare .083 ed at i amete 001. es in t	Nagelkerl R Square .1 iteration r estimates the Equation	ce e 111			
-Continue <u>Logistic Regression Screen</u> <i>Covariate box will read</i> SES(cat) (Use Save and Options buttons as needed) -Ok	Step 1 <sup>a</sup> a. Vat π(x), s Equat SES C (indica (1/.17) SES C (indica ((.37))	ses ses(1) ses(2) Constant riable(s) er success pro tion: logit( Odds Ratio ator referentimes hights $81 = 5.61 \times 1000$ Odds Ratio ator referentimes hights Odds Ratio	B 1.725 989 .956 intered b is pr $(\pi(x))$ b: odds ince cat her for K higher b: odds ince cat gher for Q k higher (Q k higher)	S.E. .253 .210 .175 on ste ob of = .95 slow/c egory low i er for smed/ egory pr med	Wald $47.358$ $46.488$ $22.152$ $29.673$ ep 1: ses.attending a $6 - 1.725$ coddshigh = $7$ is high)income thahigh incomoddshigh = $7$ is high)dium incomor high incom	$\frac{df}{2}$ $1$ $1$ $1$ $1$ $e^{-}$ $n h$ $e t$ $e^{-}$ $n e t$	Sig. .000 .000 .000 .000 .000 demic .989d 1.725 igh in han lc 989 = than h	Exp(B) $.178$ $.372$ $2.600$ $e school$ $= .1781$ $acome)$ $bw incom$ $= .3719$ $aigh incom$	ne)

Example 38:	Respor	Response: Y – (0-non-academic, 1-academic)										
ANALYZE>REGRESSION>	Explan	Explanatory:										
BINARY LOGISTIC	Gender	Gender (0-female, 1-male)										
Logistic Regression Screen	School	School (0-public, 1 – private)										
-Move y to Dependent box		Model Summary										
-Move Gender and School to	C.											
Covariate box	Step	-2 Log like	elinood	Cox & Snell R Sqt	Cox & Shell K Square							
Covariate box will read	1	792.754 <sup>a</sup> .062 .08				.083						
Gender	a. Estima	tion terminated	l at iteratio	n number 4 because pa	arameter es	timates changed by less	than					
School	001			Ĩ		6 9						
-Method: Enter	.001.			<b>T</b> 7 • 11 • /1	<b>T</b>							
-Click Categorical Button	r		1	Variables in th	e Equation	1		<b></b>				
Covariate Screen		-	В	S.E.	Wald	df	Sig.	Exp(B)				
-Move Gender and School from	Step 1 <sup>a</sup>	gender(1)	070	.16	.173		1 .678	.932				
Covariate box to categorical		/	- 1		•		•					

Covariates box	school(1)	1 533	272	31 701	1	000	216				
Categorical Covariate box will	school(1)	-1.555	.272	51.791	1	.000	.210				
read	Constant	1.413	.275	26.476	1	.000	4.109				
Gender(indicator)	a. Variable(s) entered on step 1: gender, school.										
School(indicator)	$\pi(x)$ , success prob is prob of attending academic school										
-Continue	Equation: $logit(\widehat{\boldsymbol{\pi}(\boldsymbol{x})}) = = 1.413 - 0.70$ gender - 1.533school										
Logistic Regression Screen	Gender Odds Ratio: oddsfemale/oddsmale = $e^{-0.70} = 1.0725$										
Covariate box will read	(indicator reference category male)										
Gender(cat)	(odds 1.0275 times higher to attend academic school for females than										
School (cat)	males, other vari	ables h	eld constant)								
(Use Save and Options buttons as	School Odds Ra	tio: od	dspublic/oddspri	vate=	$e^{-1.533} = 0.2159$						
needed)	(indicator referen	nce cate	egory private)								
-Ok	(odds .2159 time	s highe	r to attend an aca	ademic	e school for publi	c sch	ools				
	than private scho	ols, oth	ner variables held	l const	tant)						
	(odds 4.6321 tim	es high	er to attend an ad	cadem	ic school for priv	ate					
	schools than public schools, other variables held constant)										

<b>Example 38</b> : (2 <sup>nd</sup> suggested approach)	Respon	nse: Y –	(0-non-acade	emic,	1-academic)					
(Note: set indicator reference on	Explan	natory:								
explanatory variable appropriately)	Gende	r (0-fema	ale, 1-male)							
ANALYZE>REGRESSION>BINARY	School	l (0-publi	ic, 1 – private	e)						
LOGISTIC	Model Summary									
Logistic Regression Screen	Cox & Spoll P									
-Move y to Dependent box	a									
-Move Gender and School to Covariate	Step	-2	2 Log likelihood		Square	_	Nagelkerke R Squa	re		
box	1		792.75	54 <sup>a</sup>	.06	2		083		
Covariate box will read	a. Estimation terminated at iteration number 4 because parameter estimates changed by									
Gender	loss than 001									
School										
-Method: Enter	r		Var	riables	in the Equation	ГТ				
-Click Categorical Button			В	S.E.	Wald	df	Sig.	Exp(B)		
Covariate Screen	Step	gender(1)	070	.169	.173	1	.678	.932		
-Move Gender and School from	1a	school(1)	1 522	272	21 701	1	000	1 62 1		
Covariate box to categorical Covariates	-	sciiooi(1)	1.555	.272	51.791	1	.000	4.034		
box		Constant	120	.128	.889	1	.346	.887		
-Highlight School	a. Varia	ble(s) enter	ed on step 1: gen	der, sc	hool.					
-Change to Indicator(first)	$\pi(x), s$	uccess pi	ob is prob of	fatte	nding acaden	nic s	school			
Categorical Covariate box will read	Equat	ion: logi	$t(\widehat{\boldsymbol{\pi}(\boldsymbol{x})}) =12$	20-0	0.70gender +	1.5	33school			
Gender(indicator)	Gende	er Odds I	Ratio: oddsf	emal	e/oddsmale =	e <sup>-</sup>	0.70 = 1.0725			
School(indicator(first))	(indica	tor refer	ence categor	y ma	le)					
-Continue	(odds 1	1.0275 ti	mes higher to	o atte	nd academic	sch	ool for female	s than		
Logistic Regression Screen	males)		C							
Covariate box will read	School	l Odds F	<b>Ratio</b> : oddspr	ivate	/oddspublic					
Gender(cat)	(indica	ator refe	rence catego	ory p	ublic)					
School (cat)	School	l Odds R	atio: oddspri	vate/	oddspublic =	$e^{1.!}$	$5^{533} = 4.6321$			
(Use Save and Options buttons as	(odds 4	4.6321 ti	mes higher to	o atte	nd an acaden	nic	school for priv	ate		
needed)	school	s than pu	blic schools,	othe	r variables h	eld	constant)			
-Ok		I	,							

## **Example 5: From notes:**

Model Mo:  $\beta_0 + \beta_1 x + \beta_2$  school  $+\beta_3 d_1 + \beta_4 d_2$ where we use two dummy variables to characterize the SES variable.

1. Ho: being in an academic program is independent from SES when controlling for academic score and school type (i.e. Ho:  $\beta_3 = \beta_4 = 0$ ) versus Ha: Ho is not true,  $\alpha = 0.05$ 

2. Assumptions are met

3. Test Statistic  $\chi_o^2 = -2(L_o - L_1) = -2(-332.058 - -325.980) = -2(-332.058 + 325.980) = 12.156$ , df = 5 -3 = 2 (from following output)

4. P-valule =  $P(\chi^2 > 12.156) < 0.005$  (by table) (from following output)

5. Reject Ho since the p-value is  $< \alpha$ .

6. At a significance level of 5%, the data suggests that SES is associated with the chance of a student being in an academic program even when controlling for academic score and school type.

 $L_0 = \log$  likelihood under the Ho model  $L_1 = \log$  likelihood under the H<sub>1</sub> model

df = difference in the number of parameters of the two models

There are many ways to find  $L_0$  and  $L_1$  and the Test Statistic  $\chi_0^2$  in SPSS. We will examine several.

Characteristics of the Binary Regression approach:

- •It only allows a lowest value reference category on the response variable.
- •It allows for a switch in reference category on any explanatory variable.
- •It allows for a stepwise entry of parameters.
- •It allows for the saving of predicted probabilities and residuals.

Characteristics of the GLM approach:

- •It allows for the binary reference category to be switched on the response variable.
- •It only allows for all high value reference category (or all low value) on all explanatory variables.

•It allows for the saving of predicted probabilities, confidence intervals, and residuals. And much more!

Note that the 12.147 in the examples below is the number you get when you use  $L_0$  and  $L_1$  to more decimal places. So the 12.147 corresponds to the 12.156 from the classroom notes.

ANALYZE>REGRESSION>BINARY	*Response: Y – (0 – nonacademic, 1 – academic)
LOGISTIC	Explanatory:
Logistic Regression Screen	X – academic score (total of reading, writing, math, science, civics)
-Move y to Dependent box	School (0 – public, 1 – private)
-Move x, School to Covariate box	

# Example 5 REGRESSION PATH: MODEL L<sub>0</sub> (NO SES). METHOD ENTER

Covariate box will read	Model Summary								
X		-210	σ	Cox & Snell R		Nagelkerke F	2		
School	Ster			Carrows		Carrows	`		
-Method: Enter	Step	likelino	od	Square		Square	_		
-Click Categorical Button	1 664.107 <sup>a</sup>				.243	.3	24		
Covariate Screen	a. Estimation terminated at iteration number 4 because								
-Move School from Covariate box to	parameter estimates changed by loss than 001								
categorical Covariates box	Variable in the Exception								
Categorical Covariate box will read			\ \	ariables in	the Equ	lation			
School(indicator)			В	S.E.	Wald	df	Sig.	Exp(B)	
-Continue	Step 1 <sup>a</sup>	х	.02	7 .003	99.25	51 1	.000	1.027	
Logistic Regression Screen		school(1)	-1 42	1 288	24 33	37 1	000	241	
Covariate box will read		3011001(1)	-1.42	1 .200	24.30	1	.000	.271	
X		Constant	-5.73	8	58.95	52 1	.000	.003	
School (cat)	a. Variał	ble(s) entered	d on step	1: x, school	•				
(Use Save and Options buttons as	$\pi(\mathbf{x}), \mathbf{su}$	iccess pro	b is pro	ob of atter	<mark>iding</mark>	academic s	chool		
needed)	Equatio	on: logit( $\pi$	$(\overline{\mathbf{x}}) = -5$	.738 + 0.0	27x -1.	.421school			
-Ok	Total T	est Score	Â1p	oint increa	se in th	ne total test	score x		
	multipli	ies the odd	s of a s	student bei	ng in a	n academic	e progran	n by a	
$-2L_0 = 664.107$	factor o	of $e^{0.027} = 1$	.027, c	other varial	oles he	ld constant	. The od	ds of	
$L_0 = -332.058$	being in	n an acadei	nic pro	gram incre	ease by	2.7% for 6	each unit	increase	
	in x, oth	her variabl	es held	constant.					
	School	<b>Odds Rat</b>	io: odd	lspublic/od	ldspriv	ate = $e^{-1.4}$	$^{21} = .241$	5	
	(indicat	tor refere	nce cat	egory priv	vate)				
	(odds .2	2415 times	higher	to attend	academ	nic school f	for public	c schools	
	than pri	vate schoo	ols, oth	er variable	s held	constant)	1		
	(odds 1/	/.2415 = 4	.1413 t	imes highe	er to att	tend acader	nic schoo	ol for	
	private	schools the	an publ	lic schools	, other	variables h	eld cons	tant)	

# Example 5: REGRESSION PATH: MODEL L<sub>1</sub> (INCLUDES SES), METHOD ENTER

ANALYZE>REGRESSION>	*Response: $Y - (0 - nonacademic, 1 - academic)$										
BINARY LOGISTIC	Explanatory:	X – academic score	e (total	of reading	ng, v	vriting	, math,				
Logistic Regression Screen	science, civics)										
-Move y to Dependent box	School (0 – public, 1 – private)										
-Move x, School, SES to Covariate box	SES (1-low, 2-middle, 3-high)										
Covariate box will read	(categorize SES as 2 dummies)										
Х	$d_1 = 0 - not low, 1 - low$										
School	$d_2 = 0 - not med, 1 - medium$										
SES	Model Summary										
-Method: Enter	-		a (								
-Click Categorical Button	~		Cox 8	& Snell R							
Covariate Screen	Step	-2 Log likelihood	S	quare	N	agelkerke	e R Square				
-Move School, SES from Covariate box to	1	651.960 <sup>a</sup>		.258			.345				
categorical Covariates box	a. Estimation term	inated at iteration number 5	because p	arameter estin	nates o	changed b	by less than				
Categorical Covariate box will read	.001.										
School(indicator)		Variables	in the Eq	uation	- 1						
SES (indicator)		В	S.E.	Wald	df	Sig.	Exp(B)				

 ${}_{\rm Page}42$ 

-Continue	Sten	v	025	003	80 320	1	000	1.025	
Logistic Regression Screen	1a	л	.025	.005	00.520	1	.000	1.025	
Covariate box will read	1	school(1)	-1.346	.291	21.369	1	.000	.260	
X		ses			11.890	2	.003		
School (cat)		ses(1)	891	.285	9.756	1	.002	.410	
SES (cat)		(2)			0.070			404	
(Use Save and Options buttons as needed)		ses(2)	706	.235	8.978	1	.003	.494	
-Ok		Constant	-4.722	.801	34.798	1	.000	.009	
$-2L_1 = 651.960, L_1 = -325.980$	a. Varia	ble(s) entered on st	ep 1: x. school. ses.						
<b>SES Odds Ratio</b> : oddslow/oddshigh =	$\pi(\mathbf{x}), \mathbf{y}$	success prob	is prob of a	ttendi	ng acade	mic	schoo	l	
$e^{-0.891} = 0.4102$	Equation: $logit(\widehat{\pi(x)}) = -4.722 + 0.025x - 1.346 school - 0.891d1 - 0.706d2$								
(indicator reference category is high)	<b>Total Test Score:</b> A 1 point increase in the total test score x								
(odds 0.4102 times higher to attend	multiplies the odds of a student being in an academic program by a								
academic school for low income than high	factor	of $e^{0.025} = 1$ .	025, other va	riables	s held cor	istar	t. The	e odds of	
income, other variables held constant)	being	in an acaden	nic program i	ncreas	e by 2.5%	for	each i	unit	
$(odds \ 1/0.4102 = 2.4378 \text{ times higher to})$	increa	se in x, other	variables he	ld cons	stant.				
attend academic school for high income	Schoo	l Odds Rati	o: oddspublic	c/oddsi	orivate =	$e^{-1}$	346 = .	2603	
than low income, other variables held	(indic	ator referen	ce category	privat	<b>e</b> )	-			
constant)	(odds	.2603 times ]	higher to atte	nd aca	demic scl	iool	for pu	blic	
SES Odds Ratio: oddsmed/oddshigh =	schoo	ls than privat	e schools, oth	ner var	iables he	ld co	onstant	;)	
$e^{706} = .4936$	(odds	1/.2603 = 3.8	8417 times hi	gher to	attend a	cade	emic so	chool for	
(indicator reference category is high)	privat	e schools tha	n public scho	ols, ot	her varia	bles	held c	onstant)	
(odds 0.4936 times higher to attend	1		1	,				,	
academic school for medium income than									
high income, other variables constant)									
(odds 1/0.4102 = 2.0259 times higher to									
attend academic school for high income									
than low income, other variables constant)									

There is another approach for obtaining the necessary L<sub>o</sub> and L<sub>1</sub> for this problem. Example 5: REGRESSION PATH: ALL PARAMETERS, METHOD FORWARD (LR)

Example J. REORESSION PATH. ALL	ALL FARAMETERS, METHOD FOR WARD (LR)									
ANALYZE>REGRESSION>BINARY	*F	Response: Y – (0 –	- nonacademic,	1 – academic)						
LOGISTIC	Explanatory	: X – academic sco	ore (total of read	ling, writing, math,						
Logistic Regression Screen	science, civi	cs)								
-Move y to Dependent box	School $(0 - \text{public}, 1 - \text{private})$									
-Move x, School, and SES to Covariate		SES (1-low	v, 2-middle, 3-h	igh)						
box	(SES characterized with 2 dummy variables)									
Covariate box will read	$d_1 = 0 - not low, 1 - low$									
X	$d_2 = 0 - not med, 1 - medium$									
School	Model Summary									
SES	Step	-2 Log likelihood	Cox & Snell R	Nagelkerke R Square						
-Method: Forward (LR)			Square							
-Click Categorical Button	1	692.268 <sup>a</sup>	.207	.276						
Covariate Screen	23	651.960 <sup>b</sup>	.243	.345						
-Move School, SES from Covariate box	a. Estimation ter	rminated at iteration nur	mber 4 because							
to categorical Covariates box	parameter estim	ates changed by less that	un .001.							
Categorical Covariate box will read		Omnibus Tests of	Model Coefficients							
School(indicator)		Chi-square	df	Sig.						
SES (indicator)	Step 1 St	ep 139.082	1	.000						
-Continue	Blo	ck 139.082	1	.000						
Logistic Regression Screen				-						

 $_{\rm Page}43$ 

Covariate box will read		Model	139.082		1			000
v	Step	2 Step	28.161		1			.000
A Sahaal (aat)		Block	167.243		2			.000
School (cat)		Model	167.243		2			.000
SES (cat)	Step	3 Step	12.147		2			.002
(Use Save and Options buttons as		Block	179.390		4			.000
needed)		Model	179.390		4			.000
-Ok	Variables in the Equation							
$-2L_0 = 664.107$			D	0 F	- 	10	а.	F (D)
$L_0 = -332.058$	C to a		B 027	5.E.	Wald		S1g.	Exp(B)
	1 <sup>a</sup>	Х	.027	.003	100.411	1	.000	1.028
$-2L_1 = 651,960$		Constant	-7.055	.695	103.104	1	.000	.001
$L_1 = -325,080$	Step	х	.027	.003	99.251	1	.000	1.027
$L_1 = -525.960$	20	school(1)	-1.421	.288	24.337	1	.000	.241
	~	Constant	-5.738	.747	58.952	1	.000	.003
Step 3: Chi-square 12.147	Step	х	.025	.003	80.320	1	.000	1.025
	5-	school(1)	-1.346	.291	21.369	1	.000	.260
		ses			11.890	2	.003	
		ses(1)	891	.285	9.756	1	.002	.410
		ses(2)	706	.235	8.978	1	.003	.494
		Constant	-4.722	.801	34.798	1	.000	.009
	a. Variał	ole(s) entered o	n step 1: x.					
	b. Varial	ble(s) entered of	on step 2: schoo	1.				
	c. Varial	ole(s) entered o	n step 3: ses.					
			Mode	el if Term	Removed			
			Mode	l Log	Change in -2 Lo	g		Sig. of the
		Variable	Likeli	ihood	Likelihood		df	Change
	Step	) 1 x	-415	.675	139.082		1	.000
	Step	x x	-396	.464	128.820		1	.000
	Stor	scho	-346	.134 324	28.161		1	.000
	Siep	son son	-373	157	24 354		1	.000
		sence	-332	.053	12.147		2	.002
	π	(x), succes	s prob is p	rob of a	ttending ac	ade	mic sel	hool
	Equati	ion: logit( $\pi$	(x) = -4.722	2+0.025	x-1.346scho	01-0	.891d1	-0.706d2

# Example 5 GLM PATH: MODEL L<sub>0</sub> (NO SES) LOG LIKELIHOOD KERNAL (to obtain L<sub>0</sub>)

Analysis>Generalized Linear	*Response: Y – (0 – nonacademic, 1 – academic)						
Models>Generalized Linear Models	Explanatory:						
Type of Model Tab:	X – academic score (total	of reading,	writing, m	ath, science, civics)			
Binary Logistic	School (0 – public, 1 – private)						
Response Tab:	SES (1-low, 2-middle, 3-high) (characterize with 2 dummies)						
Dependent Variable: y	$d_1 = 0 - not low, 1 - low$						
(change Binary reference category to	$d_2 = 0 - not med, 1 - medium$						
first – lowest value)	Goodness of Fit <sup>b</sup>						
Predictor Tab:		Value	df	Value/df			
x in the Covariate box	Deviance	598.611 598.611	535	1.119			
school in Factors box	Pearson Chi-Square	533.771	535	.998			
Model Tab:	Scaled Pearson Chi-Square	533.771	535				
x as a main effect	Log Likelihood <sup>a</sup>	-332.053					
school as main effect	Akaike's Information Criterion	670.107					
$\checkmark$ intercept in model	Finite Sample Corrected AIC	670.147					
Statistics tab:	(AICC)						

Page4

Log-Likelihood Function: Kernal All other tabs: left as defaults.	Bayesian Information Criterion (BIC)       683.298         Consistent AIC (CAIC)       686.298         Dependent Variable: y Model: (Intercept), x, school a. The kernel of the log likelihood function is displayed and used in computing information criteria.         Parameter Estimates								
$-2L_0 = 664.107$ $L_0 = -332.058$	Parameter Est Parameter	timates B	Std. Error	95% Profile Lik Interval	elihood Confidence	Hypothesis Test			
				Lower	Upper	Wald Chi- Square	df	Sig.	
	(Intercept) x [school=.00] [school=1.00] (Scale)	-5.738 .027 -1.421 0 <sup>a</sup> 1 <sup>b</sup>	.7474 .0027 .2881	-7.236 .022 -2.009	-4.301 .032 875	58.952 99.251 24.337	1 1 1	.000 .000 .000	
	Dependent Variable: y Model: (Intercept), x, school a. Set to zero because this parameter is redundant. b. Information criteria are in small-is-better form. $\pi(x)$ , success prob is prob of attending academic school								
	Equation: log	$git(\pi(x))$	))=-5.	738 + 0.027x	-1.421school				

### Example 5 MODEL L<sub>1</sub>: GLM APPROACH: LOG LIKELIHOOD KERNAL (to obtain L<sub>1</sub>)

Enumple 5 MODEL EL. OEM I	III ROMON L										
Analysis>Generalized Linear											
Models>Generalized Linear	*Response: Y –	(0 - nona)	academic	z, 1 – acaden	nic)						
Models	Explanatory: X	– academ	ic score	(total of read	ling,	writing, math,	science, civics)				
Type of Model Tab	School $(0 - pub)$	lic, 1 – pi	rivate)								
Pinery Logistic	SES (1-low, 2-m	11ddle, 3-	high)	`							
Billary Logistic	(View SES as an d = 0 mot low	1 X With	h 3 levels	5)							
Response Tab:	$d_1 = 0 - \text{not low}$	= 0 - 100100, $1 - 100$									
Dependent Variable: y	$u_2 = 0 = not met$	$_2 = 0 - \text{not med}, 1 - \text{medium}$									
(change Binary reference	Г — — — — — — — — — — — — — — — — — — —		Goodile	Value	df	Value/df					
category to first – lowest	Deviance			632,552	571	1 1(	08				
value)	Scaled Deviance			632.552	571	1.10					
Predictor Tab:	Pearson Chi-Squ	are		570.021	571	.99	98				
x in the Covariate box	Scaled Pearson C	Chi-Square	e	570.021	571						
school in Easters box	Log Likelihood <sup>a</sup>			-325.980							
	Akaike's Informa	ation Crite	rion	661.960							
ses in Factors box	Finite Sample Co	orrected A	IC	662.061							
Model Tab:	(AICC)		-								
x as a main effect	Bayesian Inform	ation Crite	erion	683.944							
school as main effect	(BIC)	a LTA)		600.044							
ses in Factors box	Consistent AIC (	CAIC)		688.944		-					
✓ intercept in model	Model: (Intercent	ble: y	ol ses								
Statistics tab:	a. The kernel of t	he log like	elihood fu	nction is disp	laved	and used in					
Log-Likelihood Function:	computing inforr	nation crit	eria.	I I I I I I I I I I I I I I I I I I I							
Kornal	b. Information cr	iteria are i	n small-is	-better form.							
		Test	ts of Mod	el Effects							
All other tabs:				Type III							
left as defaults.	Source	Likeliho	ood Ratio	df		Sig					
	(Intercent)	CIII-,		a	1	51g.					
$-2L_1 = 651.960$	(intercept)		70.08 98.68	2	1	.000					
$L_1 = -325.980$	school		24.354	4	1	.000					
	ses	ses 12.147 2 .002									
ses Likelihood Ratio Chi-	Dependent Varia	ble: y									
Square	Model: (Intercept	t), x, schoo	ol, ses	. —							
12 1/7	│┏━━━━━		Pa	arameter Est	imate	S					
12.17/	Deremeter	р	Std.	95% Pro	tile Li	kelihood	Hypothesis Test				
	r drameter	D	CITOL	Conna	ence I	merval	rypomesis rest				

 $_{\rm Page}45$ 

			Lower	Upper	Wald Chi- Square	df	Sig.	
(Intercept)	-4.722	.8005	-6.319	-3.176	34.798	1	.000	
Х	.025	.0028	.020	.031	80.320	1	.000	
[school=.00]	-1.346	.2912	-1.939	793	21.369	1	.000	
[school=1.00]	0 <sup>a</sup>							
[ses=1.00]	891	.2853	-1.455	335	9.756	1	.002	
[ses=2.00]	706	.2355	-1.173	248	8.978	1	.003	
[ses=3.00]	0 <sup>a</sup>							
(Scale)	1 <sup>b</sup>							
Dependent Variable: y Model: (Intercept), x, school, ses a. Set to zero because this parameter is redundant. b. Fixed at the displayed value.								
$\pi(\mathbf{x})$ , success	prob is	prob	of attending	g academic so	chool			
Equation: log	$it(\widehat{\pi(x)})$	=-4.72	2+0.025x-1.	346school-0.	891d1-	0.70	)6d2	

# Model Building and Application in Logistic Regression

# **Choosing Predictors:**

-At least 10 of each outcome for every predictor

-Avoid too many predictors

-Avoid multi-collinearity (caused by linear dependence of predictors)

-If have interaction terms in model, all lower order terms should be included

-all of none of a dummy variable should be included

-Include all variables important to design of study, even if effect non-significant.

HSB model (could have up to 29 predictors (because 292/10 = 29)

Analyze > Descriptive Statistics > Frequencies Variable(s): y

Ok

	у												
				Valid	Cumulative								
		Frequency	Percent	Percent	Percent								
Valid	.00	292	48.7	48.7	48.7								
	1.00	308	51.3	51.3	100.0								
	Total	600	100.0	100.0									

# Finding Models:

Stepwise selection:

 $_{\rm Page}46$ 

Backward: Drop existing variable with the "least significant' effect at each step. Stop when another step shows no further improvement in the model

Forward: Add variable that shows the "most significant" effect at each step. Stop when another step shows no further improvement in the model.

Example: HSB with predictors x (= reading + writing + math + science + civics), science, SES and school (Model M2 for our purposes here)

Y – (0-non-academic, 1-academic)	(	Categorical	l Variab	les Cod	lings						
X – academic score (total of reading,				Daram	ater coding						
writing, math, science, civics)		Fre	mency	(1)	(2)	5					
School (0-public, 1 – private)	ses	1.00	139	1.00	000	)					
Science (score on Science test)	565	1.00	107	11000							
SES (1-low, 2-middle, 3-high)		2.00	299	.000	0 1.000	)					
(categorize SES with 2 dummies)	school	3.00	162 506	000. 000	000.						
$d_1 = 0 - not low, 1 - low$	senoor	.00	200	.00							
$d_2 = 0 - not med, 1 - medium$		1.00	94	1.000	0						
ANALYZE>REGRESSION>BINARY		0	mnibus	Tests o	f Model C	Coefi	ficients	5			
LOGISTIC				Chi-sq	uare		df	,	Sig.		
Logistic Regression Screen	Step 1	Step		1	39.082			1	.000		
-Move y to Dependent box		Block		1	39.082			1	.000		
-Move SES, school, x, science to		Model		1	39.082			1	.000		
Covariate box	Step 2	Step		1	28.161			1	.000		
Covariate box will read		Model		1	67.243			2	.000		
Х	Step 3	Step		-	13.506			1	.000		
School	_	Block		1	80.749			3	.000		
science		Model		1	80.749			3	.000		
SES	Step 4	Step			12.806			2	.002		
-Method: Forward LR		Block		1	93.555			5	.000		
-Click Categorical Button		Model			93.555			5	.000		
Covariate Screen				Vari	ables in th	ie E	quatio	n			
-Move School and SES from Covariate									95% C.I.f	or EXP(B)	
box to categorical Covariates box			В	S.E.	Wald	df	Sig.	Exp(B)	Lower	Upper	
-Highlight School	Step 1 <sup>a</sup>	х	.027	.003	106.411	1	.000	1.028	1.022	1.033	
-Change to Indicator(first)	G. Oh	Constant	-7.055	.695	103.104	1	.000	.001	0.055	7.007	
Categorical Covariate box will read	Step 2 <sup>o</sup>	school(1)	1.421	.288	24.337	1	.000	4.143	2.355	1.033	
School(indicator(first))		Constant	-7.160	.714	100.606	1	.000	.001	11022	11000	
SES(indicator)	Step 3 <sup>c</sup>	school(1)	1.449	.291	24.699	1	.000	4.257	2.404	7.537	
-Continue		X	.040	.005	72.679 12.879	1	.000	1.041 030	1.031	1.051	
Logistic Regression Screen		Constant	-7.308	.725	101.481	1	.000	.001	.907	.712	
Covariate box will read	Step 4 <sup>d</sup>	ses			12.506	2	.002				
X		ses(1)	926	.289	10.256	1	.001	.396	.225	.698	
school (cat)		ses(2)	731	.238	9.413	1	.002	.481	.302	.768	
science		school(1)	1.374	.295	21.738	1	.000	3.952	2.218	7.043	
ses(Cat)		x	.039	.005	65.371	1	.000	1.039	1.030	1.049	
Ok		science	065	.018	13.502	1	.000	.937	.905	.970	
		Constant	-6.193	.786	62.014	1	.000	.002			

 $P_{age}4$ 

Options: Check CI for exp(B) 95% Continue Ok	a. Variable( b. Variable( c. Variable( d. Variable(	(s) entered on step 1: x. (s) entered on step 2: school. (s) entered on step 3: science (s) entered on step 4: ses.	2. Model Summary					
AIC = $-2(\log \text{ likelihood } - \# \text{ parameters})$	Step	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square				
$= -2\log \text{ likelihood} + 2(\text{#parameters})$	1 2 3	692.268 <sup>a</sup> 664.107 <sup>b</sup> 650.601 <sup>b</sup>	.207 .243 .260	.276 .324 .347				
M1: (predictors x, school, and science) AIC = $650.601 + 8 = 658.601$ (Parameters for intercept, x, school, and science)	<ul> <li>4 637.795<sup>b</sup> .276 .3</li> <li>a. Estimation terminated at iteration number 4 because parameter estimates changed by less than .001.</li> <li>b. Estimation terminated at iteration number 5 because parameter estimates changed by less than .001</li> </ul>							
M2: (predictors x, school, science, and SES) AIC = 637.795 + 12 = 649.795 (Parameters for intercept, x, school, science, and 2 dummy variables for SES) AIC SMALLER – BETTER MODEL	<b>π(x) (sue</b> <b>Eqn</b> : log 0.926d <sub>1</sub>	<b>ccess prob is prob</b> git(π(x))= = -6.193 - 0.731d <sub>2</sub>	<b>of attending academ</b> 3 + 0.039x+1.374schoo	<b>tic school)</b> ol -0.065science-				

# **Classification Tables:**

If  $\widehat{\pi(x)}$  is >= 0.5, predict success for this individual more likely than failure (Classify Individual as Success) If  $\widehat{\pi(x)}$  is < 0.5, predict failure for this individual less likely than success (Classify Individual as Failure)

## Create Classification table to compare actual measurements with predicted categories.

Model M2 above provides the following classification table. We are interested in the one at Step 4, when all predictors are in the model.

Our model has percentage corrects of over 70% for both Y=0 and Y=1.

Classification Table<sup>a</sup>

-						
-	-		Predicte	d		
			/	Percentage		
	Observed	.00	1.00	Correct		
Step 1	y .00	198	94	67.8		
	1.00	81	227	73.7		
	Overall Percentage			70.8		
Step 2	у .00	208	84	71.2		
	1.00	73	235	76.3		
	Overall Percentage			73.8		
Step 3	у .00	202	90	69.2		
	1.00	75	233	75.6		

	Overall Percentage			72.5
Step 4	у .00	208	84	71.2
	1.00	74	234	76.0
	Overall Percentage			73.7

a. The cut value is .500

If you run M2 with Enter, and put all the predictors in at once, you obtain:

		Classifica	tion Table <sup>a</sup>					
	_		Predicted					
			S					
	Observed		.00	1.00	Percentage Correct			
Step 1	у	.00	208	84	71.2			
		1.00	74	234	76.0			
	Overall P	ercentage			73.7			

a. The cut value is .500

# **Correlation:**

Run model M2 as above, entering all variables at once, but pull up the Save pop-up, and choose Probabilities under predicted values when you run. This will create a column labeled PRE\_1 in your data file.

Then follow the path Analyze > Correlate > Bivariate and put the variables y and Predicted Probability [PRE\_1] in the Variables box. Say ok. You will receive the output found below.

(I suspect that Karen ran a GLM approach, which will also provide predicted values, given that her correlation is slightly different from mine, viz. 0.533 versus 0.536.)

Correlations							
			Predicted				
		у	probability				
у	Pearson	1	.536**				
	Correlation						
	Sig. (2-tailed)		.000				
	Ν	600	600				
Predicted	Pearson	.536**	1				
probability	Correlation						
	Sig. (2-tailed)	.000					
	Ν	600	600				

\*\*. Correlation is significant at the 0.01 level (2-tailed).

 $R^2 = .536^2 = 0.287$ 

When you run model M2 as above, and enter all the variables at once, the Model Summary returned looks like this. Note that it reminds you that it did one step in order to obtain the output.

Model Summary									
-2 Log Cox & Snell R Nagelkerke R									
Step	likelihood	Square	Square						
1	637.795 <sup>a</sup>	.276	.368						

a. Estimation terminated at iteration number 5 because parameter estimates changed by less than .001.

# "Influence Diagnostics"

How influential single observations are on the model fit is of interest. Some measures that look at this include.

1. DFBETA – the change in the parameter estimate divided by its standard error, when excluding the individual 2. c- confidence displacement – this measures the change in the joint confidence intervals for all parameters when excluding the individual

3. The change in the Pearson  $\chi^2$  or  $G^2$  (Deviance) when the observation is omitted

The larger the measure, the greater the influence of the individual.

We can run model M2 as above, bring up the save box and select Leverage values and DFBeta(s) under the Influence group. Say Continue and Ok to run.

The leverage values in SPSS consider observations to be influential when they have values larger than 2(k+1)/nwhere n is the sample size, and k is the number of predictors.

To create a descriptive statistics box to look at a summary of your results, follow the path Analyze>Descriptive Statistics > Descriptives and bring the Leverage value variable, and the DFBETA variables to the Variable(s) box. Drag and drop them to get them to appear in the same order as in Karen's notes.

Descriptive Statistics										
		Minimu	Maximu		Std.					
	Ν	m	m	Mean	Deviation					
Leverage value	600	.00255	.03216	.0100000	.00502988					
DFBETA for x	600	00116	.00050	.0000000	.00019665					
DFBETA for	600	00250	.00346	.0000000	.00072190					
science	600	07060	05180	0000021	01207167					
school(1)	000	07009	.05189	0000021	.01207107					

# **Descriptive Statistics**

DFBETA for ses(1)	600	04132	.03338	.0000014	.01171271
DFBETA for ses(2)	600	03367	.03515	0000003	.00967199
Valid N (listwise)	600				

We can see that there are no standout individual values by looking at these summary statistics.

# **Multicategory Logit Models – Nominal Response**

Variable Y has J categories, i = 1, ..., J

 $\pi_1, \ldots, \pi_J$  are the probabilities that observations fall into the categories

Question: What effect do certain predictors have on the  $\pi_1, \ldots, \pi_J$ ?

Binary Logistic: Model odds for success (probability of success in relationship to probability of failure) Multicategory Logit: Simultaneously model all relationships between probabilities for pairs of categories(model odds of falling within one category instead of another)

Baseline Logits approach: Pair a baseline category with all remaining categories Usually last category (J)

Define: Baseline Logit  $Log(\pi_i/\pi_J), i = 1,... J-1$ 

Model with one predictor x  $Log(\pi_i/\pi_J) = \alpha_i + \beta_i x, i = 1, ..., J-1$ 

For each category i compared, a new set of parameters introduced

SPSS models these J-1 model equations simultaneously.

Baseline category logit model allows comparison of categories a and b

 $Log(\pi_a/\pi_b) = log((\pi_a/\pi_J)/(\pi_b/\pi_J)) = log(\pi_a/\pi_J) - log(\pi_b/\pi_J) = (\alpha_a - \alpha_b) - (\beta_a - \beta_b)x$ 

Example: Contraceptive Use in Dependency on Age File is Contraceptive.sav 3165 currently married women Variables contra (s(sterilization), o(other), n(none)) age (midpoints of 5 year intervals, viz: 17.5, 22.5...47.5) a2 – age squared

Other columns:

freq -frequency in each age midpoint/contraception level cell s – frequencies in s category for each age midpoint o – frequencies in o category for each age midpoint n - frequencies in n category for each age midpoint  $lsn - log (\pi_s/\pi_n)$  $lon - log (\pi_o/\pi_n)$ 

It is of interest to create a crosstab table, for visual acuity. We will also find a Chi-square statistic, in order to test for independence. (Below, we see that we reject our null hypothesis of no relationship between age and method of contraception, and conclude that the two variables are related (associated) )

Commands: Data>Weight Cases Weight cases by Frequency Variable: freq OK

Analyze >Descriptive Statistics > Crosstabs Row(s): age Column(s): contra Statistics button: check chi square Continue OK

# age \* contra Crosstabulation

Count

			contra					
		n	n o s					
age	17.50	232	61	3	296			
	<mark>22.50</mark>	<mark>400</mark>	<mark>137</mark>	<mark>80</mark>	<mark>617</mark>			
	27.50	301	131	216	648			
	32.50	203	76	268	547			
	37.50	188	50	197	435			
	42.50	164	24	150	338			
	47.50	183	10	91	284			
Total		1671	489	1005	3165			

The o and n numbers were entered into their own columns in the contraceptive.sav file.

# **Chi-Square Tests**

Page 52

	Value	df	Asymp. Sig.
	value	ui	(2-slueu)
Pearson Chi-	430.028 <sup>a</sup>	<mark>12</mark>	.000
Square			
Likelihood Ratio	<mark>521.103</mark>	<mark>12</mark>	.000
N of Valid Cases	3165		

a. 0 cells (.0%) have expected count less than 5. The minimum expected count is 43.88.

We will be looking at using a multicategory logit model.

It is therefore of interest to graph age versus  $\log(\pi_s/\pi_n)$  and age versus  $\log(\pi_o/\pi_n)$ . These log values can be obtained for our age midpoints by using the models below, but they can also be readily obtained if we note that for age midpoints of interest:  $\log(\pi_s/\pi_n) = \log(s/n)$  and  $\log(\pi_o/\pi_n) = \log(o/n)$ 

Commands: Transform>Compute Variable Target Variable lsn Numeric Expression: LN(s/n)

Transform>Compute Variable Target Variable:lon Numeric Expression: LN(o/n)

Graph>Legacy Dialogs>Scatter/Dot> Overlay Scatter Define X-Y pairs Pair Y Variable X Variable 1 lsn age 2 lon age Ok

Double click to activate graph

Click on lsn on legend, and change type to a triangle in the popup, apply, close

Right click on an lsn triangle (or a lon circle), select add fit line at total from the dropdown menu, select quadratic in the popup, and apply.

Close the activated graph window to return it to the output file.



**SES Odds Ratio**: oddslow/oddshigh =  $e^{-0.891} = 0.4102$ 

### (reference category is high)

(odds 0.4102 times higher to attend academic school for low income than high income, other variables held constant)

(odds 1/0.4102 = 2.4378 times higher to attend academic school for high income than low income, other variables held constant)

**SES Odds Ratio:** oddsmed/oddshigh =  $e^{-.706} = .4936$ 

### (reference category is high)

(odds 0.4936 times higher to attend academic school for medium income than high income, other variables held constant)

(odds 1/0.4102 = 2.0259 times higher to attend academic school for high income than low income, other variables held constant) We see that a quadratic model appears to fit our data well.

We are interested in comparing two models, a linear model and a quadratic model. We will obtain logit equations, and perform tests of significance. Our results are summarized here, and the output and commands needed to obtain the output are provided below.

Linear Model (baseline n)	Quadratic Model (baseline n)
$Log(\pi_s/\pi_n) = \alpha_s + \beta_{1s}age$	$Log(\pi_s/\pi_n) = \alpha_s + \beta_{1s}age + \beta_{2s}age^2$
$Log(\pi_o/\pi_n) = \alpha_o + \beta_{1o}age$	$Log(\pi_0/\pi_n) = \alpha_0 + \beta_{10}age + \beta_{20}age^2$
From output BELOW:	From output BELOW:
$Log(\widehat{\pi s}/\widehat{\pi n}) = -2.236 + 0.053age$	$Log(\hat{\pi s}/\hat{\pi n}) = -12.618 + .710 age010 age^2$
$Log(\hat{\pi o}/\hat{\pi n}) =152 - 0.037$ age	$Log(\hat{\pi o}/\hat{\pi n}) = -4.550 + .264age005age^2$

Linear Model (baseline o)	Quadratic Model (baseline o)	$\overline{4}$
From output BELOW:	From output BELOW:	പ്ര
$Log(\widehat{\pi s}/\widehat{\pi o}) = -2.084 + 0.091 age$	$Log(\hat{\pi s}/\hat{\pi o}) = -8.068 + .446age - 0.005age^2$	ag

EXAMPLE OF INTERPRETATION: Given that the chosen contraceptive is either sterilization or none, find how the odds that a woman used sterilization at 26 changed from the odds that a woman used sterilization at 25.

Log $(\widehat{\pi s}/\widehat{\pi n}) = -12.6 + .710$ age -.010age<sup>2</sup>  $(\widehat{\pi s}/\widehat{\pi n}) = e^{(-12.6 + .710$ age -.010age<sup>2</sup>)} For age = 25,  $(\widehat{\pi s}/\widehat{\pi n}) =$ odds (25) =  $e^{(-12.6 + .710(25) -.010(25)2)}$ For age = 26,  $(\widehat{\pi s}/\widehat{\pi n}) =$ odds (26) =  $e^{(-12.6 + .710(26) -.010(26)2)}$ Odds(26)/Odds(25) =  $e^{(.710(26-25) - .010(262-252))} = e^{(0.710(1) - .010(51))} = 1.22$ There was a 22% increase in the odds from the age of 25 to 26.

Linear Model (Test of Fit):	Quadratic Model (Test of Fit):
From output BELOW	From output BELOW
Pearson $\chi^2 = 268.169$ with df = 10, p-value =0.000	Pearson $\chi^2 = 18.869$ with df = 8, p-value =0.016
Deviance = 293.979 with df =10, p-value = 0.000	Deviance = $20.475$ with df = 9, p-value = $0.009$

The quadratic model shows a marginally "acceptable" model.

Linear Model (Wald $\chi^2$ signif test on odds ratios)					Quadrati	Quadratic Model(Wald $\chi^2$ signif test on odds ratios)				
OUTPU	OUTPUT BELOW					OUTPUT BELOW				
Odds-Variable $\chi^2$ DfP-valueOdds-Variable $\chi^2$ DfP-value							P-value			
ratio					ratio					
Ster-	age	130.295	1	<.001	Ster-	age	232.22	1	<.001	
None					None					
Other-	age	33.008	1	<.001		age <sup>2</sup>	218.25	1	<.001	
None					Other-	age	31.47	1	<.001	
Ster-	age	168.869	1	<.001	None	_				
Other						age <sup>2</sup>	39.23	1	<.001	
					Ster-	age	54.79	1	<.001	
					Other	_				
						age <sup>2</sup>	28.24	1	<.001	

Here we see that both age and  $age^2$  have a significant effect on the log odds for sterilization versus none, other versus none, and sterilization versus other.

# **QUADRATIC MODEL COMMANDS AND OUTPUT:**

To obtain parameter estimate output for 1<sup>st</sup> two equations:

Model button:

Levels (categories) for the variable contra will be n, o, s (and the first will be the one lowest in the alphabet) Run with baseline logit of n (none) <u>Commands:</u> Analyze>Regression>Multinomial Logistic Dependent: contra Reference Category: first Covariate(s): age a2 Multinomial Logistic Regression: Model popup window

Choose Custom

Build Terms: Main effects

Select age and a2 from Factors & Covariates window and bring them over to Forced Entry Terms window. Continue

Statistics Window: leave defaults (Under model: Case processing summary, Pseudo R-square, Step summary, Model fitting information, and under Parameters: Estimates and Likelihood ratio tests, and under Define Subpopulations: Covariate patterns defined by factors and covariates) Also check "Goodness of Fit" under model.

Criteria Window: leave defaults

Options Window: leave defaults

Save Window: select estimated response probabilities

<u>To obtain parameter output estimate for 3<sup>rd</sup> equation</u>, rerun above commands using baseline logit other. When you click on reference category, it offers a custom option. Set it to o (other).

	Parameter Estimates									
			C. 1					95% Confidence Interval fo		
	ļ	( I	Std.	1		1 1		Елр	( <b>D</b> )	
contra	i <sup>a</sup>	В	Error	Wald	df	Sig.	Exp(B)	Lower Bound	Upper Bound	
0	Intercept	<mark>-4.550</mark>	.694	42.998	1	.000				
	age	<mark>.264</mark>	.047	<mark>31.472</mark>	1	.000	1.302	1.187	1.428	
	a2	<mark>005</mark>	.001	<mark>39.233</mark>	1	.000	.995	.994	.997	
s	Intercept	<mark>-12.618</mark>	.757	277.545	1	.000				
	age	<mark>.710</mark>	.046	243.225	1	.000	2.033	1.860	2.223	
	a2	<mark>010</mark>	.001	218.250	1	.000	.990	.989	.992	

a. The reference category is: n.

# **Parameter Estimates**

								95% Confiden	ce Interval for
			Std.					Exp	(B)
contr	a <sup>a</sup>	В	Error	Wald	df	Sig.	Exp(B)	Lower Bound	Upper Bound
n	Intercept	4.550	.694	42.998	1	.000			
	age	264	.047	31.472	1	.000	.768	.700	.842
	a2	.005	.001	39.233	1	.000	1.005	1.003	1.006
s	Intercept	<mark>-8.068</mark>	.949	72.285	1	.000			
	age	<mark>.446</mark>	.060	<mark>54.789</mark>	1	.000	1.561	1.388	1.757
	a2	<mark>005</mark>	.001	28.843	1	.000	.995	.993	.997

a. The reference category is: o.

The Goodness-of-Fit output for this quadratic model can be found on both runs.

Goodness-of-Fit				
Chi-Square	df	Sig.		

- -----

Pearson	18.869	8	.016
Deviance	20.475	8	.009

# LINEAR MODEL COMMANDS AND OUTPUT

The Linear model commands are similar to the QUADRATIC MODEL COMMANDS. The only difference is that there is only one explanatory variable, age, rather than two explanatory variables, age and a2.

# **Parameter Estimates**

								95% Confiden	ce Interval for
			Std.					Exp	(B)
con	ntra <sup>a</sup>	В	Error	Wald	df	Sig.	Exp(B)	Lower Bound	Upper Bound
0	Intercept	<mark>152</mark>	.190	.638	1	.424			
	age	<mark>037</mark>	.006	33.008	1	.000	.964	.951	.976
s	Intercept	-2.236	.159	198.226	1	.000			
	age	.053	.005	130.295	1	.000	1.055	1.045	1.065

a. The reference category is: n.

### **Parameter Estimates**

								95% Confiden	ce Interval for
			Std.					Exp	(B)
cont	ra <sup>a</sup>	В	Error	Wald	df	Sig.	Exp(B)	Lower Bound	Upper Bound
n	Intercept	.152	.190	.638	1	.424			
	age	.037	.006	33.008	1	.000	1.038	1.025	1.051
s	Intercept	<mark>-2.084</mark>	.215	<mark>93.818</mark>	1	.000			
	age	<mark>.091</mark>	.007	168.869	1	.000	1.095	1.080	1.110

a. The reference category is: o.

**Goodness-of-Fit** 

	Chi-Square	df	Sig.	
Pearson	268.169	10	.000	
Deviance	293.978	10	.000	

# TO TEST IF AGE HAS AN EFFECT ON CONTRACEPTIVES:

WE WILL NEED **<u>FURTHER QUADRATIC MODEL OUTPUT</u>** HERE, VIZ: A TABLE WITH MODEL FITTING INFORMATION, AS BELOW. IT IS OBTAINED IN THE OUTPUT FROM THE COMMANDS ABOVE.

Ho: all the age and age<sup>2</sup> parameters are 0 ( $\beta_{1s} = \beta_{2s} = \beta_{1o} = \beta_{2o} = 0$ ), and the intercept model suffices) versus Ha: at least one of the parameters is non-zero, choose  $\alpha = 0.05$ Assumptions hold

 $\chi^2 = -2(L_o - L_1) = 500.628$  with df =4 (see output below)

Where  $L_0$  is intercept only model and  $L_1$  is model with age and age<sup>2</sup>

P-value < 0.001

Reject Ho at the 5% significance level. The data provides sufficient evidence that age has an effect on the choice of contraceptive.

	Model Fitt	ting Inform	nation		
	Model				
	Fitting				
	Criteria	Likeli	Likelihood Ratio Tests		
	-2 Log				
	Likeliho	Chi-			
Model	od	Square	df	Sig.	
Intercept	601.377				
Only					
Final	100.749	<mark>500.628</mark>	<mark>4</mark>	.000	

# COMPARING ESTIMATED PROBABILITIES FROM CONTINGENCY TABLES WITH PROBABILITIES ESTIMATED FROM THE QUADRATIC MODEL

Karen's notes have the formulas for calculating the  $\hat{\pi}$  s.

# Commands:

SPSS calculates the  $\hat{\pi}$  s for us when we run the quadratic model as we did above, and, in addition, use the Save button, and in the Saved variables box, check "estimated response probabilities".

This will yield three new columns in the data file contraceptive.sav. They will be labeled, EST1-1, EST2\_1, AND EST3\_1. For age= 22.5, we have

EST1 1 EST2 2 EST3 3 .64 .23 .13

And these values match the values found in Karen's notes for  $\hat{\pi n}$ ,  $\hat{\pi o}$  and  $\hat{\pi s}$  when age = 22.5

# **Ordinal Response**

Model

Y is an ordinal variable with categories 1, 2,...J Cumulative probability for category  $j = P(Y \le J) = \pi_1 + ... + \pi_j$ , j = 1,2,...J $P(Y \le 1) \le P(Y \le 2) \le ... P(Y \le J) = 1$ 

Model an ordinal response:

Model the cumulative response probabilities (or cumulative odds)

Cumulative Logits:

 $\text{Logit} \left( \mathsf{P}(\mathsf{Y}{<=}j) \right) = \log(\frac{P(Y{\leq}j)}{1{-}P(Y{\leq}j)}) = \log\big(\frac{\pi_1{+}\pi_2{+}\cdots{+}\pi_j}{\pi_{j{+}1}{+}\cdots{+}\pi_j}\big)$ 

Categories are then: I  $(1 \ 2 \ \dots \ j)$  and II  $(j+1 \ \dots \ J)$ 

For one predictor variable x, proportional odds model is , for j = 1, 2, ... J-1,

Logit  $(p(Y \le j)) = \alpha_j + \beta x$ 

NOTE: slope  $\beta$  <u>the same</u> for all cumulative logits

 $\beta$  describes the effect of x on the odds of falling into categories 1 to j.

The model assumes that this effect is the same for all cumulative odds.

NOTE: Property of this model (see textbook page 181 for diagram) is that, for  $\beta >0$ , as x increases, the probability of falling in a lower category increases.

NOTE: SPSS models logit(Y<=j) =  $\alpha_j - \beta x$ And reports the negative of the slope!

Example: HSB Data: Y, the ordinal variable here, is SES (low – 1, middle – 2, high – 3) Math Writing Sex (1 – Male, 2 - Female) Race (1-Hispanic 2-Asian 3 - Black 4- White)

We will use the model equation below for j = 1, 2. Logit(P(Y<=j) =  $\alpha j+\beta_1$ math + $\beta_2$ writing+ $\beta_3$ sex +  $\beta_{41}$ race1 to  $\beta_{42}$ race2 +  $\beta_{43}$  race3

We first check to be sure we have no ses/sex cells or ses/race cells with very few observations, as this would cause problems in analyzing our model.

<u>Commands:</u> Analyze>Descriptive Statistics > Crosstabs Row(s): ses Column(s): sex Ok

Analyze>Descriptive Statistics > Crosstabs Row(s): ses Column(s): race Ok

# ses \* sex Crosstabulation

Count

_		se	sex		
		1.00	2.00	Total	
ses	1.00	50	89	139	
	2.00	144	155	299	
	3.00	79	83	162	
Total		273	327	600	

# ses \* race Crosstabulation

Count

		race						
	1.00	2.00	3.00	4.00	Total			
1.00	23	8	24	84	139			
2.00	34	14	24	227	299			
3.00	14	12	10	126	162			
	71	34	58	437	600			
	1.00 2.00 3.00	1.00           1.00         23           2.00         34           3.00         14           71	ra           1.00         2.00           1.00         23         8           2.00         34         14           3.00         14         12           71         34	race1.002.003.001.00238242.003414243.00141210713458	race1.002.003.004.001.0023824842.003414242273.00141210126713458437			

There is no indication of empty or "small" cells to worry about.

We will create our estimated model. Some commands below are included to help us check our fit.

COMMANDS: Analyze > Regression >Ordinal Dependent: ses Factor(s): sex race Covariate(s): writing math Output button: Display area: (Leave defaults: Goodness of fit statistics, Summary statistics, parameter estimates) Check "Test of parallel lines" Saved Variables area: Check "Predicted category" Continue Ok

Output of use from these commands is posted below with explanations.

		8	-	
	-2 Log		-	
Model	Likelihood	Chi-Square	df	Sig.
Intercept	1221.075			
Only				
Final	1148.593	<mark>72.482</mark>	<mark>6</mark>	<mark>.000</mark>

# **Model Fitting Information**

Link function: Logit.

Here we see that the model that includes all the predictors fits significantly better than the model that omits all the predictors (the intercept model): Reject Ho (intercept only model), with  $\chi^2 = 74.482$ , df =6, and p-value <0.001.

# Goodness-of-Fit

	Chi-Square	df	Sig.
Pearson	1148.222	1138	.410
Deviance	1125.438	1138	.599

Link function: Logit.

Here both the Pearson  $\chi^2$  (1148.22, df = 1138) and the Deviance (1125.438, df = 1138) suggest good fit.

# Pseudo R-Square

Cox and	.114
Snell	
Nagelkerke	.130
McFadden	.058

Link function: Logit.

The Pseudo  $R^2$  values give the impression that the model fits somewhat well.

Here are the parameter estimates.

**Parameter Estimates** 

							95% Confidence Interva	
		Estimate	Std. Error	Wald	df	Sig.	Lower Bound	Upper Bound
Threshold [ses = 1.00]		2.509	.555	20.422	1	.000	1.421	3.597

	[ses = 2.00]	4.927	.587	70.497	1	.000	3.777	6.077
Location	writing	.027	.011	5.833	1	.016	.005	.049
	math	.043	.012	14.121	1	.000	.021	.066
	[sex=1.00]	.456	.170	7.210	1	.007	.123	.788
	[sex=2.00]	0 <sup>a</sup>			0			
	[race=1.00]	143	.255	.314	1	.575	644	.357
	[race=2.00]	151	.345	.193	1	.660	827	.524
	[race=3.00]	476	.279	2.910	1	.088	-1.024	.071
	[race=4.00]	0 <sup>a</sup>		•	0	•	•	

Link function: Logit.

a. This parameter is set to zero because it is redundant.

Threshold values. The estimates tell us the values used to differentiate socioeconomic classes. 2.509 differentiates low ses from middle ses and 4.927 differentiates middle ses from high ses.

# Math estimate:

Math is a significant predictor of SES ( $\chi^2 = 14.121$ , df = 1, p-value < 0.001)

A one point increase in math results in an increase in the ordered log odds of being in a higher ses category of 0.043 (all other variables in the model held constant). Since  $e^{.043} = 1.044$ , a one point increase in math increases your chance of being in a higher SES category by 4.4%.

# Writing estimate:

Writing is a significant predictor of SES ( $\chi^2 = 5.833$ , df = 1, p-value < 0.016)

A one point increase in writing results in an increase in the ordered log odds of being in a higher ses category of 0.027 (all other variables in the model held constant). Since  $e^{.027} = 1.027$ , a one point increase in math increases your chance of being in a higher SES category by 2.7%.

# Sex estimate:

Sex is a significant predictor of SES ( $\chi^2 = 7.210$ , df = 1, p-value < 0.007)

If a person is male, this results in increase in the ordered log odds of being in a higher socioeconomic category of 0.456 (all other variables held constant). Since  $c^{.456} = 1.578$ , males are 57.8% more likely to be in a higher SES category than females (all other variables held constant). ALTERNATIVE INTERPRETATION. Is it better to say that the ordered log odds of being in a higher socioeconomic category are 1.578X times higher for males than females (all other variables held constant)?

Race estimate:

Race is not a significant predictor of SES. (all p-values for all race dummy variables are large)

We need to check if it is reasonable to assume that the slopes for the cumulative odds are all the same (that is: that the S-curves for all SES categories are truly parallel. We are checking if the chance of falling into a higher SES is equally affected by the predictor variables for all SES categories.)

	rest of raraner Lines										
	-2 Log										
Model	Likelihood	Chi-Square	df	Sig.							
Null	1148.593										
Hypothesis											
General	1142.142	6.450	6	.375							

Tost of Dorollol Linosa

The null hypothesis states that the location parameters (slope coefficients) are the same across response categories. a. Link function: Logit.

Here Ho is that the lines are parallel. Our result is non-significant, indicating that we cannot reject the possibility that the effect is the same for all categories. This supports our model, but we should, as always be cautious when "accepting" our Ho because we don't know the error probability for that decision.

Commands above have created a Predicted Response Category Variable in the HSB data file. It will be to the far right of all the columns in your HSB data file should be labeled PRE\_# (# will depend on whether this is your first predictor variable for this set of data, or if you have been doing several examinations are keeping the new columns that are appearing in the HSB data file when you do), and when you hover over the label, you will see that it is a column of Predicted Response Category.

We create a crosstab table of SES and the predicted response categories for SES with our model in order to see how well the model does.

# **COMMANDS**

Analyze>Descriptive Statistics>Crosstabs Row(s): ses Column(s): Predicted Response Category OK

# ses \* Predicted Response Category Crosstabulation

		Predicted			
		1.00	2.00	3.00	Total
ses	1.00	18	117	4	139
	2.00	10	259	30	299
	3.00	7	131	24	162
Total		35	507	58	600

Because all observations do not fall on the diagonal, we know that our model is not very powerful for prediction of the SES of an individual.



# **Loglinear Models for Contingency Tables**

Consider the GLM:

Link function g(mu) = log(mu) Poisson -random component Linear predictor - systematic component

Contingency table: Response Variable: Count Predictors: Categorical Variables that define the contingency table

Suppose have two categorical variables, X (with I levels) and Y(with J levels)

 $\pi_{ij}$  = joint probability =  $\pi_{i+}\pi_{+j}$ ,  $i = 1, \dots, I$  and  $j = 1, \dots, J$ 

Expected cell count  $\mu_{ij} = n\pi_{ij}$ , i = 1, ..., I and j = 1, ..., J

If the two variables are independent, Expected cell count  $\mu_{ij} = n\pi_{i+}\pi_{+j}$  i = 1, ..., I and j = 1, ..., J

 $Log(\mu ij) = log(n) + log(\pi_{i+}) + log(\pi_{+j}) i = 1,..., I and j = 1,... J$ 

This is the **loglinear model of independence.** It only fits if the two variables are independent.

Writing convention:

 $Log(\mu_{ij}) = \lambda + \lambda_i^X + \lambda_j^Y i = 1, ..., I \text{ and } j = 1, ... J$ 

-  $\lambda$  is the overall effect (included to ensure the  $\Sigma \pi_{ij} = n$ )

-  $\lambda_i^X$  is the (main or marginal) effect of category i of variable X on the log of the expected cell count

-  $\lambda_j^{Y}$  is the (main or marginal) effect of category j of variable Y on the log of the expected cell count

Model fit statistics for the model (Pearson's  $\chi^2$  and the Likelihood Ratio Statistic) are used to test if the two variables are independent.

Example: Two way table: (alien\_walle.sav)

 Movie\_n (1 – Alien, 2-Wall-e) Rating (12 – 1&2, 34 – 3&4, 56 - 5&6, 78 – 7&8, 910 – 9&10)

We will examine the fit of the model mentioned above with this data.

Commands: Analyze>Loglinear>General Distribution of Cell Counts: Poisson Factor(s): Movie\_n Rating Model tab: Custom Build Term(s): Choose Main effects Highlight movie\_n and rating in the Factors and Covariates box Move them (with the arrow) to the Terms in Model box Continue Save tab (just note there are no default options and cancel) Options tab (just note the defaults and cancel) Display defaults: Frequencies, Residuals Plot defaults: Adjusted residuals, Normal probability for adjusted Continue OK

NOTE: At no point did we tell it a dependent count variable here. It looked at the data file and made its own cell counts and residuals table. It matches the movie-n values and ratings values with the counts from the freq column. Pretty cool, eh. HOWEVER,

It is probably better to do a Data>Weight Cases command, though, as it might not always be the case that it is obvious in what column SPSS should look for the frequencies.

(See how it counted)

### Cell Counts and Residuals<sup>a,b</sup>

	-	Obse	erved	Expected			Standardized	Adjusted	
movie_n	rating	Count	%	Count	%	Residual	Residual	Residual	Deviance
1.00	12	1932	.6%	3938.206	1.2%	-2006.206	-31.969	-43.810	-35.505
	34	1760	.5%	2022.127	.6%	-262.127	-5.829	-7.933	-5.962
	56	8403	2.5%	8101.193	2.4%	301.807	3.353	4.662	3.333
	78	54233	16.4%	46835.648	14.1%	7397.352	34.181	55.721	33.336
	910	83874	25.3%	89304.827	26.9%	-5430.827	-18.173	-38.596	-18.362
2.00	12	6758	2.0%	4751.794	1.4%	2006.206	29.104	43.810	27.349
	34	2702	.8%	2439.876	.7%	262.124	5.307	7.931	5.216
	56	9473	2.9%	9774.807	2.9%	-301.807	-3.053	-4.662	-3.069

 $P_{age}65$ 

,	78	49114	14.8%	56511.352	17.1%	-7397.352	-31.118	-55.721	-31.837
	910	113185	34.2%	107754.17	32.5%	5430.827	16.544	38.596	16.408
				3					

a. Model: Poisson

b. Design: Constant +  $movie_n + rating$ 

You will get a lot more output, but we are only interested in the Goodness of Fit tests table.

Goodness-of-Fit Tests<sup>a,b</sup>

	Value	df	Sig.
Likelihood Ratio	4823.100	4	.000
Pearson Chi-	4692.375	4	.000
Square			

a. Model: Poisson

b. Design: Constant + movie\_n + rating

Ho: independence between movie\_n and rating versus Ha: is a relationship (dependence) between movie\_n and rating

Reject Ho (Pearson  $\chi^2$ =4692.4, df = 4, p-value <0.001). We do not have independence.

# WE SHOULD NOT USE THIS MODEL. WE NEED A MODEL THAT TAKES INTERACTION INTO EFFECT. WE DO NOT HAVE INDEPENDENCE OF RATING AND MOVIE.....

However, for interest ONLY, and to help us get ready to read the output for more complex models, I include the parameter estimates here.

					95% Confidence Interval					
		Std.			Lower	Upper				
Parameter	Estimate	Error	Z	Sig.	Bound	Bound				
Constant	11.588	.003	4210.148	.000	11.582	11.593				
[movie_n =	188	.003	-53.821	.000	195	181				
1.00]										
[movie_n =	0 <sup>a</sup>	•	•		•	•				
2.00]										
[rating = 12]	-3.121	.011	-284.761	.000	-3.143	-3.100				
[rating = 34]	-3.788	.015	-250.354	.000	-3.818	-3.758				
[rating = 56]	-2.400	.008	-307.255	.000	-2.415	-2.385				
[rating = 78]	645	.004	-168.046	.000	653	638				
[rating = 910]	0 <sup>a</sup>	•	•							

Parameter Estimates<sup>b,c</sup>

a. This parameter is set to zero because it is redundant.

b. Model: Poisson

c. Design: Constant + movie\_n + rating

# FOR INTEREST ONLY, WE LOOK AT READING THIS OUTPUT, AND INTERPRETING IT IF IT WERE REASONABLE TO USE IT.

X-rating i = 1, 2, 3, 4, 5 Y-movie j = 1, 2

$$\begin{split} & \text{Log}(\mu_{11}) = \lambda + \lambda_1^{\text{Rating}} + \lambda_1^{\text{Movie}} = 11.588 - .3.121 - .188 \\ & \text{Log}(\mu_{21}) = \lambda + \lambda_2^{\text{Rating}} + \lambda_1^{\text{Movie}} = 11.588 - 3.788 - .188 \\ & \text{Log}(\mu_{31}) = \lambda + \lambda_3^{\text{Rating}} + \lambda_1^{\text{Movie}} = 11.588 - .2.400 - .188 \\ & \text{Log}(\mu_{41}) = \lambda + \lambda_4^{\text{Rating}} + \lambda_1^{\text{Movie}} = 11.588 - .645 - .188 \\ & \text{Log}(\mu_{51}) = \lambda + \lambda_5^{\text{Rating}} + \lambda_1^{\text{Movie}} = 11.588 - .188 \\ & \text{Log}(\mu_{12}) = \lambda + \lambda_1^{\text{Rating}} + \lambda_2^{\text{Movie}} = 11.588 - .188 \\ & \text{Log}(\mu_{12}) = \lambda + \lambda_1^{\text{Rating}} + \lambda_2^{\text{Movie}} = 11.588 - .188 \\ & \text{Log}(\mu_{22}) = \lambda + \lambda_3^{\text{Rating}} + \lambda_2^{\text{Movie}} = 11.588 - 3.788 \\ & \text{Log}(\mu_{22}) = \lambda + \lambda_3^{\text{Rating}} + \lambda_2^{\text{Movie}} = 11.588 - 3.788 \\ & \text{Log}(\mu_{32}) = \lambda + \lambda_3^{\text{Rating}} + \lambda_2^{\text{Movie}} = 11.588 - .645 \\ & \text{Log}(\mu_{42}) = \lambda + \lambda_5^{\text{Rating}} + \lambda_2^{\text{Movie}} = 11.588 - .645 \\ & \text{Log}(\mu_{52}) = \lambda + \lambda_5^{\text{Rating}} + \lambda_2^{\text{Movie}} = 11.588 - .645 \end{split}$$

Examples of interpretation: (IF OUR MODEL WERE TRUE AND WE HAD INDEPENDENCE)

 $Log(\mu_{i1}/\mu_{i2}) = \lambda_1^{Movie} - \lambda_2^{Movie} = \lambda_1^Y - \lambda_2^Y = -0.188$ does not depend on the level i of X(rating). That is, it is the same for all levels I of X(rating).

 $e^{-.188} = 0.8286 = \text{oddsAlien/oddsWalle}$  is the same for each rating

The odds of a movie being Alien are 0.8286 as high as the odds of a movie being Walle, for each rating. The odds of a movie being Walle are 1/.8286 = 1.2068 times higher than the odds of a movie being Alien, for each rating.

 $Log(\mu_{1j}/\mu_{5j}) = \lambda_1^{Rating} - \lambda_5^{Rating} = \lambda_1^X - \lambda_5^X = -3.121$ does not depend on the level j of Y(movie). It is the same for all levels, j, of Y. That is, it is the same for both movies.

 $e^{-3.121} = .0441 =$ odds Rating 1/odds Rating 5 is the same for both movies.

The odds of a rating of 1 are .0441 times as high as the odds of a rating of 5, for each movie. The odds of a rating of 5 are 1/0.441 = 22.6757 times higher than a rating of 1, for each movie.

# TWO WAY SATURATED MODEL (using movie\_n, rating, and movie\_n\*rating)

Model is:  $Log(\mu ij) = \lambda + \lambda_i^{Movie} + \lambda_j^{Rating} + \lambda_{ij}^{MovieRating}$  j = 1,2 and i = 1,2,3,4,5

$$\begin{split} &\text{Movie\_n (1 - Alien, 2-Wall-e)} \\ &\text{Rating (12 - 1\&2, 34 - 3\&4, 56 - 5\&6, 78 - 7\&8, 910 - 9\&10)} \\ &\text{Log}(\mu_{11}) = \lambda + \lambda_1^{\text{Rating}} + \lambda_1^{\text{Movie}} + \lambda_{11}^{\text{Movie*Rating}} \end{split}$$

$Log(\mu_{21}) = \lambda + \lambda_2^{Rating}$	$+  \lambda_1{}^{Movie} + +  \lambda_{12}{}^{Movie*Rating}$
$Log(\mu_{31}) = \lambda + \lambda_3^{Rating}$	$+ \lambda_1^{Movie} + + \lambda_{13}^{Movie*Rating}$
$Log(\mu_{41}) = \lambda + \lambda_4^{Rating}$	$+ \lambda_1^{Movie} + + \lambda_{14}^{Movie*Rating}$
$Log(\mu_{51}) = \lambda + \lambda_5^{Rating}$	$+  \lambda_1{}^{Movie} + \lambda_{15}{}^{Movie*Rating}$
$Log(\mu_{21}) = \lambda + \lambda_1^{Rating}$	$+  \lambda_2{}^{Movie} + \lambda_{21}{}^{Movie*Rating}$
$Log(\mu_{22}) = \lambda + \lambda_2^{Rating}$	$+ \lambda_2^{Movie} + \lambda_{22}^{Movie*Rating}$
$Log(\mu_{32}) = \lambda + \lambda_3^{Rating}$	$+ \lambda_2^{Movie} + \lambda_{23}^{Movie*Rating}$
$Log(\mu_{42}) = \lambda + \lambda_4^{Rating}$	$+ \lambda_2^{Movie} + \lambda_{24}^{Movie*Rating}$
$Log(\mu_{52}) = \lambda + \lambda_5^{Rating}$	$+ \ \lambda_2{}^{Movie} + \ \lambda_{25}{}^{Movie*Rating}$

Note: a saturated model gives a perfect fit for any sample. Since we seek a smaller set of parameters in order to have a more parsimonious model to describe the population, a saturated model is not our goal. Commands below are included in order that students can run this model, and see that the residuals are all zero when the saturated model is fit.

**COMMANDS:** Data> Weight Cases Weight cases by: Frequency Variable: freq Ok Analyze>Loglinear>General Distribution of Cell Counts: Poisson Factor(s): Rating Movie n Model tab: Choose Saturated, Continue Save tab (just note there are no default options and cancel) Options tab (just note the defaults and cancel) Display defaults: Frequencies, Residuals Plot defaults: Adjusted residuals, Normal probability for adjusted OK

		Observ	/ed	Expect	Expected		Standardized	Adjusted	
movie_n	rating	Count	%	Count	%	Residual	Residual	Residual	Deviance
1.00	12	1932.500	.6%	1932.500	.6%	.000	.000	.000	.000
	34	1760.500	.5%	1760.500	.5%	.000	.000	.000	.000
	56	8403.500	2.5%	8403.500	2.5%	.000	.000	.000	.000
	78	54233.500	16.4%	54233.500	16.4%	.000	.000	.000	.000
	910	83874.500	25.3%	83874.500	25.3%	.000	.000	.000	.000
2.00	12	6758.500	2.0%	6758.500	2.0%	.000	.000	.000	.000
	34	2702.500	.8%	2702.500	.8%	.000	.000	.000	.000
	56	9473.500	2.9%	9473.500	2.9%	.000	.000		.000
	78	49114.500	14.8%	49114.500	14.8%	.000	.000		.000
	910	113185.500	34.1%	113185.500	34.1%	.000	.000	.000	.000

SPSS adds .5 to each observed count for the saturated model. Later, I will show you how to change this. Cell Counts and Residuals<sup>a,b</sup>

### a. Model: Poisson

b. Design: Constant + movie\_n + rating + movie\_n \* rating Note that for the saturated model, observed and expected counts are exactly the same.

### THREE WAY TABLE MODEL

# THREE WAY EXAMPLE: ALIEN WALL-E

For each gender, we create a cross tab tables for movie \* rating.

Commands Analysis>Descriptive Statistics >Crosstabs Row(s): movie Column(s): rating Layer: sex

Count

### movie \* rating \* sex Crosstabulation

-				rating							
sex			12	34	56	78	910	Total			
f	movie	Alien	430	356	1313	5075	6777	13951			
		Wall-e	1104	441	1274	5998	19106	27923			
	Total		1534	797	2587	11073	25883	41874			
m	movie	Alien	1502	1404	7090	49158	77097	136251			
		Wall-e	5654	2261	8199	43116	94079	153309			
	Total		7156	3665	15289	92274	171176	289560			

We will examine various models for looking at this data when we take the categorical predictors movie, gender and rating all into account.

# MODEL 1: SATURATED (3 WAY INTERACTION, ALL 2 WAY INTERACTION, ALL MAIN EFFECTS)

 $Log(\mu_{ijk}) = \lambda + \lambda_i^{movie} + \lambda_j^{rate} + \lambda_k^{sex} + \lambda_{jk}^{movierate} + \lambda_{ik}^{moviesex} + \lambda_{ij}^{sexrate} + \lambda_{ijk}^{movieratesex}$ I = 1,2, j = 1,...5, k = 1,2

Commands: Analyze>Loglinear>General Distribution of Cell Counts: Poisson Factor(s): Movie\_n Rating Sex\_n Model tab: Saturated Continue Save tab (just note there are no default options and cancel) Options tab (just note the defaults and cancel) Display defaults: Frequencies, Residuals Plot defaults: Adjusted residuals, Normal probability for adjusted OK

Goodness-of-Fit Tests<sup>a,b</sup>

	Value	df	Sig.
Likelihood Ratio	.000	0	
Pearson Chi-	.000	0	
Square			

a. Model: Poisson

b. Design: Constant + movie\_n + rating + sex\_n +

movie\_n \* rating + movie\_n \* sex\_n + rating \*

sex\_n + movie\_n \* rating \* sex\_n

	-	_	Obser	ved	Expected			Standardized	Adjusted	
movie_n	rating	sex_n	Count	%	Count	%	Residual	Residual	Residual	Deviance
1.00	12	1.00	430.500	.1%	430.500	.1%	.000	.000	.000	.000
		2.00	1502.500	.5%	1502.500	.5%	.000	.000	.000	.000
	34	1.00	356.500	.1%	356.500	.1%	.000	.000	.000	.000
		2.00	1404.500	.4%	1404.500	.4%	.000	.000	.000	.000
	56	1.00	1313.500	.4%	1313.500	.4%	.000	.000	.000	.000
		2.00	7090.500	2.1%	7090.500	2.1%	.000	.000		.000
	78	1.00	5075.500	1.5%	5075.500	1.5%	.000	.000	.000	.000
	_	2.00	49158.500	14.8%	49158.500	14.8%	.000	.000		.000
	910	1.00	6777.500	2.0%	6777.500	2.0%	.000	.000	.000	.000
		2.00	77097.500	23.3%	77097.500	23.3%	.000	.000	.000	.000
2.00	12	1.00	1104.500	.3%	1104.500	.3%	.000	.000	.000	.000
		2.00	5654.500	1.7%	5654.500	1.7%	.000	.000	.000	.000
	34	1.00	441.500	.1%	441.500	.1%	.000	.000	.000	.000
		2.00	2261.500	.7%	2261.500	.7%	.000	.000		.000
	56	1.00	1274.500	.4%	1274.500	.4%	.000	.000	.000	.000
		2.00	8199.500	2.5%	8199.500	2.5%	.000	.000		.000
	78	1.00	5998.500	1.8%	5998.500	1.8%	.000	.000	.000	.000
		2.00	43116.500	13.0%	43116.500	13.0%	.000	.000	.000	.000
	910	1.00	19106.500	5.8%	19106.500	5.8%	.000	.000		.000
		2.00	94079.500	28.4%	94079.500	28.4%	.000	.000	.000	.000

Cell Counts and Residuals<sup>a,b</sup>

a. Model: Poisson

b. Design: Constant + movie\_n + rating + sex\_n + movie\_n \* rating + movie\_n \* sex\_n + rating \* sex\_n + movie\_n \* rating \* sex\_n

					95% Confidence Interval	
Parameter	Estimate	Std. Error	Z	Sig.	Lower Bound	Upper Bound
Constant	11.452	.003	3512.569	.000	11.446	11.458
$[movie_n = 1.00]$	199	.005	-40.978	.000	209	190
$[movie_n = 2.00]$	$0^{a}$					
[rating = 12]	-2.812	.014	-205.348	.000	-2.839	-2.785
[rating = 34]	-3.728	.021	-175.198	.000	-3.770	-3.686
[rating = 56]	-2.440	.012	-211.909	.000	-2.463	-2.417
[rating = 78]	780	.006	-134.160	.000	792	769
[rating = 910]	$0^{a}$					
$[sex_n = 1.00]$	-1.594	.008	-200.891	.000	-1.610	-1.579
[sex n = 2.00]	$0^{a}$					
[movie $n = 1.00$ ] * [rating = 12]	-1.126	.029	-38.271	.000	-1.184	-1.069
[movie $n = 1.00$ ] * [rating = 34]	277	.034	-8.079	.000	345	210
[movie $n = 1.00$ ] * [rating = 56]	.054	.017	3.175	.001	.021	.087
[movie $n = 1.00$ ] * [rating = 78]	.330	.008	40.301	.000	.314	.346
[movie $n = 1.00$ ] * [rating = 910]	O <sup>a</sup>					
[movie $n = 2.00$ ] * [rating = 12]	0ª		-			
[movie_n = 2.00] * [rating = 34]	$O^a$	•	•	•	•	•
[movie $n = 2.00$ ] * [rating = 56]	$O^a$					
[movie $n = 2.001 * [rating - 78]$	O <sup>a</sup>					•
[movie_n = 2.00] [rating = 70] [movie_n = 2.00] * [rating = 910]	O <sup>a</sup>					
[movie $n = 1.00$ ] [rating = $710$ ] [movie $n = 1.00$ ] * [sev $n = 1.00$ ]	- 837	015	-56.012	000	- 867	- 808
[movie_n = $1.00$ ] [sex_n = $1.00$ ] [movie_n = $1.00$ ] * [sev_n = $2.00$ ]	037 0ª	.015	-50.012	.000	007	008
$[movie_n - 2.00] * [sev_n - 1.00]$	O <sup>a</sup>					•
$[movie_n = 2.00] + [sex_n = 1.00]$ [movie_n = 2.00] * [sex_n = 2.00]		•			·	•
[10000 - 1 - 2.00] + [500 - 1 - 2.00] [rating - 12] * [sev n - 1.00]	- 039		_1 151	250	- 105	027
$[rating = 12] * [sex_n = 1.00]$	039 O <sup>a</sup>	.034	-1.151	.230	105	.027
$[rating = 12] * [sex_n = 2.00]$	030	053	750	153	1/3	064
$[rating = 34] * [sex_n = 2.00]$	0 <i>5</i> 7	.055	750	55	145	.004
$[\text{rating} = 54] * [\text{sex}_{11} = 2.00]$	267	021	0 500		. 228	206
$[rating = 56] * [sex_n = 2.00]$	207 O <sup>a</sup>	.031	-0.300	.000	328	200
$[rating = 30] * [sex_11 = 2.00]$	278	016			400	. 247
$[rating = 78] * [sex_n = 2.00]$	378 O <sup>a</sup>	.010	-23.789	.000	409	347
$[rating = 78] + [sex_n = 2.00]$	0	•	•	•	•	•
$[rating = 910] * [sex_n = 1.00]$	0-	•		•		•
$[\text{rating} = 910] + [\text{sex_II} = 2.00]$	1 220	066	19 625		. 1.002	. 1.240
$[\text{Inovie}_{II} = 1.00] + [\text{rating} = 12] + [\text{rating} = 12] + [\text{rating} = 1.00]$	1.220	.000	18.023	.000	1.092	1.549
$[sex_n = 1.00]$	03					
$[movie_n = 1.00] * [rating = 12] *$	0"	-				
$[sex_n = 2.00]$	1 100	000	12 (00	000	0.42	1.057
$[movie_n = 1.00] * [rating = 34] *$	1.100	.080	13.698	.000	.942	1.257
$[sex_n = 1.00]$	03					
$[movie_n = 1.00] * [rating = 34] *$	$0^{a}$	•		•	·	•
$[sex_n = 2.00]$	1.010	o 1 -	22.15-		~~ ·	
$[movie_n = 1.00] * [rating = 56] *$	1.013	.045	22.466	.000	.924	1.101
$[sex_n = 1.00]$	00					
$[movie_n = 1.00] * [rating = 56] *$	$0^{a}$	•			•	•
[sex_n = 2.00]						
$[movie_n = 1.00] * [rating = 78] *$	.539	.025	21.466	.000	.490	.588
$[sex_n = 1.00]$	~					
$[movie_n = 1.00] * [rating = 78] *$	$0^{a}$	•				
[sex_n = 2.00]						
$[movie_n = 1.00] * [rating = 910]$	$0^{a}$	•	•		•	•
* [sex_n = 1.00]						
$[movie_n = 1.00] * [rating = 910]$	$0^{a}$	•				
* $[sex_n = 2.00]$	_					
$[movie_n = 2.00] * [rating = 12] *$	$0^{a}$	•				
$[sex_n = 1.00]$						
$[movie_n = 2.00] * [rating = 12] *$	$0^{a}$	•				•
$[sex_n = 2.00]$						
$[movie_n = 2.00] * [rating = 34] *$	$0^{a}$	•				•
$[sex_n = 1.00]$						
[movie_n = 2.00] * [rating = 34] *	$0^{a}$					
$[sex_n = 2.00]$						
[movie $n = 2.001 * [rating = 56] *$	$0^{a}$					
$[10000 \_ 1 = 2.00]  [1000 \_ 50]$						

 $_{\rm Page}71$ 

[movie_n = 2.00] * [rating = 56] *	0 <sup>a</sup>			
$[sex_n = 2.00]$				
[movie_n = 2.00] * [rating = 78] *	$0^{\mathrm{a}}$			
$[sex_n = 1.00]$				
[movie_n = 2.00] * [rating = 78] *	0ª			
$[sex_n = 2.00]$				
[movie_n = 2.00] * [rating = 910]	$0^{\mathrm{a}}$			
* [sex_n = 1.00]				
[movie_n = 2.00] * [rating = 910]	$0^{\mathrm{a}}$			
* [sex_n = 2.00]				

a. This parameter is set to zero because it is redundant.

b. Model: Poisson

 $c. \ Design: Constant + movie\_n + rating + sex\_n + movie\_n * rating + movie\_n * sex\_n + rating * sex\_n + movie\_n * * sex\_n + sex\_n * sex\_n *$ 

MODEL 2: HOMOGENOUS MODEL: Drop the 3 way interaction term, but include all 2 way interactions and main effects

$$\begin{split} & Log(\mu_{ijk}) = \lambda + \lambda_i^{movie} + \lambda_j^{rate} + \lambda_k^{sex} + \lambda_{jk}^{movierate} + \lambda_{ik}^{moviesex} + \lambda_{ij}^{sexrate} \\ & I = 1, 2, j = 1, \dots 5, k = 1, 2 \end{split}$$

<u>Commands</u>: Data> Weight Cases Weight cases by: Frequency Variable: freq Ok

Analyze > Loglinear >General General Loglinear Analysis Box: Distribution of Cell Counts: Poisson Factors: movie\_n rating sex\_n Model Button

Custom Main effects dropdown: Highlight factors one at a time, and click arrow to move them over to terms in model box Movie\_n Rating Sex\_n Interaction: Highlight movie\_n and rating, click arrow to bring over to terms in model box Highlight movie\_n and sex\_n, click arrow to bring over to terms in model box
Highlight rating and sex\_n, click arrow to bring over to terms in the model box Continue

Options Display: Frequencies (default) Residuals (default) Estimates (check) Plot: Adjusted Residuals (default) Normal Probability for adjusted (default) Save: (leave defaults) OK

## Goodness-of-Fit Tests<sup>a,b</sup>

	Value	df	Sig.
Likelihood Ratio	1121.297	4	.000
Pearson Chi-	1162.251	4	.000
Square			

a. Model: Poisson

b. Design: Constant + movie\_n + rating + sex\_n + movie\_n \* rating + movie\_n \* sex\_n + rating \* sex\_n

Cell Counts and Residuals<sup>a,b</sup>

-	-	_	Obse	rved	Expec	ted		Standardized	Adjusted	
movie_n	rating	sex_n	Count	%	Count	%	Residual	Residual	Residual	Deviance
1.00	12	1.00	430	.1%	232.873	.1%	197.127	12.918	15.219	11.540
		2.00	1502	.5%	1699.127	.5%	-197.127	-4.782	-15.219	-4.880
	34	1.00	356	.1%	232.127	.1%	123.873	8.130	10.602	7.532
		2.00	1404	.4%	1527.873	.5%	-123.873	-3.169	-10.602	-3.213
	56	1.00	1313	.4%	919.188	.3%	393.812	12.989	17.963	12.196
		2.00	7090	2.1%	7483.812	2.3%	-393.812	-4.552	-17.963	-4.593
	78	1.00	5075	1.5%	4455.475	1.3%	619.525	9.281	15.090	9.078
		2.00	49158	14.8%	49777.525	15.0%	-619.525	-2.777	-15.090	-2.783
	910	1.00	6777	2.0%	8111.336	2.4%	-1334.336	-14.816	-30.277	-15.252
		2.00	77097	23.3%	75762.664	22.9%	1334.336	4.848	30.277	4.834
2.00	12	1.00	1104	.3%	1301.127	.4%	-197.127	-5.465	-15.219	-5.612
		2.00	5654	1.7%	5456.873	1.6%	197.127	2.669	15.219	2.653
	34	1.00	441	.1%	564.873	.2%	-123.873	-5.212	-10.601	-5.422
		2.00	2261	.7%	2137.127	.6%	123.873	2.680	10.602	2.654

56	1.00	1274	.4%	1667.812	.5%	-393.812	-9.643	-17.963	-10.066
	2.00	8199	2.5%	7805.188	2.4%	393.812	4.458	17.963	4.421
78	1.00	5998	1.8%	6617.525	2.0%	-619.525	-7.616	-15.090	-7.739
	2.00	43116	13.0%	42496.475	12.8%	619.525	3.005	15.090	2.998
910	1.00	19106	5.8%	17771.664	5.4%	1334.336	10.009	30.277	9.888
	2.00	94079	28.4%	95413.336	28.8%	-1334.336	-4.320	-30.277	-4.330

b. Design: Constant + movie\_n + rating + sex\_n + movie\_n \* rating + movie\_n \* sex\_n + rating \* sex\_n

					95% Confide	ence Interval
Parameter	Estimate	Std. Error	Z	Sig.	Lower Bound	Upper Bound
Constant	11.466	.003	3578.338	.000	11.460	11.472
[movie_n = 1.00]	231	.005	-48.520	.000	240	221
$[movie_n = 2.00]$	$0^{a}$					
[rating = 12]	-2.861	.014	-208.969	.000	-2.888	-2.835
[rating = 34]	-3.799	.021	-179.580	.000	-3.840	-3.757
[rating = 56]	-2.503	.011	-219.850	.000	-2.526	-2.481
[rating = 78]	809	.006	-142.636	.000	820	798
[rating = 910]	0 <sup>a</sup>					
$[sex_n = 1.00]$	-1.681	.008	-220.493	.000	-1.696	-1.666
$[sex_n = 2.00]$	$0^{a}$					
[movie_n = 1.00] * [rating = 12]	936	.026	-35.622	.000	988	885
[movie_n = 1.00] * [rating = 34]	105	.031	-3.373	.001	166	044
[movie_n = 1.00] * [rating = 56]	.189	.016	11.985	.000	.158	.219
[movie_n = 1.00] * [rating = 78]	.389	.008	50.179	.000	.374	.404
[movie_n = 1.00] * [rating = 910]	$0^{a}$					
[movie_n = 2.00] * [rating = 12]	$0^{a}$					
[movie_n = 2.00] * [rating = 34]	0 <sup>a</sup>					
[movie_n = 2.00] * [rating = 56]	0 <sup>a</sup>					
[movie_n = 2.00] * [rating = 78]	$0^{a}$					
[movie_n = 2.00] * [rating = 910]	$0^{a}$					
$[movie_n = 1.00] * [sex_n = 1.00]$	554	.011	-49.926	.000	575	532
$[movie_n = 1.00] * [sex_n = 2.00]$	$0^{a}$					
$[movie_n = 2.00] * [sex_n = 1.00]$	$0^{a}$					
$[movie_n = 2.00] * [sex_n = 2.00]$	$0^{a}$					
[rating = 12] * [sex_n = 1.00]	.247	.029	8.497	.000	.190	.304
[rating = 12] * [sex_n = 2.00]	$0^{a}$					
$[rating = 34] * [sex_n = 1.00]$	.350	.040	8.785	.000	.272	.428
$[rating = 34] * [sex_n = 2.00]$	$0^{a}$					
$[rating = 56] * [sex_n = 1.00]$	.137	.022	6.135	.000	.093	.181
$[rating = 56] * [sex_n = 2.00]$	$0^{a}$					
[rating = 78] * [sex_n = 1.00]	179	.012	-14.738	.000	203	155
[rating = 78] * [sex_n = 2.00]	0 <sup>a</sup>					
[rating = 910] * [sex_n = 1.00]	$0^{a}$					
[rating = 910] * [sex_n = 2.00]	0 <sup>a</sup>					

Parameter Estimates<sup>b,c</sup>

a. This parameter is set to zero because it is redundant.

b. Model: Poisson

c. Design: Constant + movie\_n + rating + sex\_n + movie\_n \* rating + movie\_n \* sex\_n + rating \* sex\_n + rati

#### WE WISH TO TEST IF THE RATING, SEX, AND MOVIE ARE HOMOGENOUS.

Ho:  $\lambda_{ijk}^{ratesexmovie} = 0$  (homogenous) versus Ha: saturated model,  $\alpha = 0.05$ 

Reject Ho. (Pearson  $\chi^2 = 1121.297$ , df =4, pvalue <.001, Likelihood Ration = 1162.251, df =4, p-value<0.001). We need the three way interaction term. The relationship between sex, rating and movie is not homogenous.

Model 3: ALL MAIN EFFECTS, INTERACTION TERMS FOR RATE\*MOVIE AND SEX\*MOVIE (RATING\*SEX IS EXCLUDED)

Model is: Log  $\mu ijk = \lambda + \lambda_i^{sex} + \lambda_j^{rate} + \lambda_k^{movie} + \lambda_{ij}^{ratemovie} + \lambda_{ik}^{sexmovie}$  $I = 1, 2, j = 1, \dots 5, k = 1, 2$ 

Commands: Data> Weight Cases Weight cases by: Frequency Variable: freq Ok

Analyze > Loglinear >General General Loglinear Analysis Box: Distribution of Cell Counts: Poisson Factors: movie\_n rating sex\_n

Model Button Custom Main effects dropdown: Highlight factors one at a time, and click arrow to move them over to terms in model box Movie\_n Rating Sex\_n Interaction: Highlight movie\_n and rating, click arrow to bring over to terms in model box Highlight movie\_n and sex\_n, click arrow to bring over to terms in model box Continue Options Display: Frequencies (default) Residuals (default) Estimates (check)

Plot:

PIOL:

Adjusted Residuals (default)

Normal Probability for adjusted (default)

Save: (leave defaults)

OK

# Goodness-of-Fit Tests<sup>a,b</sup>

Value df Sig.			
Ũ	Value	df	Sig.

Likelihood Ratio	1599.934	8	.000
Pearson Chi-	1788.937	8	.000
Square			

b. Design: Constant + movie\_n + rating + sex\_n +

movie\_n \* rating + movie\_n \* sex\_n

$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$											
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$				Obse	rved	Expec	cted		Standardized	Adjusted	
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	movie_n	rating	sex_n	Count	%	Count	%	Residual	Residual	Residual	Deviance
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	1.00	12	1.00	430	.1%	179.447	.1%	250.553	18.704	19.766	15.826
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$			2.00	1502	.5%	1752.553	.5%	-250.553	-5.985	-19.762	-6.137
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		34	1.00	356	.1%	163.472	.0%	192.528	15.058	15.904	13.003
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$			2.00	1404	.4%	1596.528	.5%	-192.528	-4.818	-15.898	-4.921
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$		56	1.00	1313	.4%	780.484	.2%	532.516	19.061	20.598	17.346
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$			2.00	7090	2.1%	7622.516	2.3%	-532.516	-6.099	-20.598	-6.173
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $		78	1.00	5075	1.5%	5037.247	1.5%	37.753	.532	.699	.531
$\begin{array}{c c c c c c c c c c c c c c c c c c c $			2.00	49158	14.8%	49195.753	14.8%	-37.753	170	699	170
$\begin{array}{c c c c c c c c c c c c c c c c c c c $		910	1.00	6777	2.0%	7790.350	2.4%	-1013.350	-11.481	-18.140	-11.744
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$			2.00	77097	23.3%	76083.650	23.0%	1013.350	3.674	18.140	3.666
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	2.00	12	1.00	1104	.3%	1041.227	.3%	62.773	1.945	2.156	1.926
34  1.00  441  .1%  416.306  .1%  24.694  1.210  1.326  1.199    2.00  2261  .7%  2285.694  .7%  -24.694 517  -1.326 517    56  1.00  1274  .4%  1459.536  .4%  -185.536  -4.856  -5.424  -4.965    2.00  8199  2.5%  8013.464  2.4%  185.536  2.073  5.424  2.065    78  1.00  5998  1.8%  7567.153  2.3%  -1569.153  -18.038  -22.970  -18.723    2.00  43116  13.0%  41546.847  12.5%  1569.153  7.698  22.970  7.651    910  1.00  19106  5.8%  17438.779  5.3%  1667.221  12.625  22.402  12.432    2.00  94079  28.4%  95746.221  28.9%  -1667.221  -5.388  -22.402  -5.404			2.00	5654	1.7%	5716.773	1.7%	-62.773	830	-2.156	832
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$		34	1.00	441	.1%	416.306	.1%	24.694	1.210	1.326	1.199
56  1.00  1274  .4%  1459.536  .4%  -185.536  -4.856  -5.424  -4.965    2.00  8199  2.5%  8013.464  2.4%  185.536  2.073  5.424  2.065    78  1.00  5998  1.8%  7567.153  2.3%  -1569.153  -18.038  -22.970  -18.723    2.00  43116  13.0%  41546.847  12.5%  1569.153  7.698  22.970  7.651    910  1.00  19106  5.8%  17438.779  5.3%  1667.221  12.625  22.402  12.432    2.00  94079  28.4%  95746.221  28.9%  -1667.221  -5.388  -22.402  -5.404			2.00	2261	.7%	2285.694	.7%	-24.694	517	-1.326	517
2.00  8199  2.5%  8013.464  2.4%  185.536  2.073  5.424  2.065    78  1.00  5998  1.8%  7567.153  2.3%  -1569.153  -18.038  -22.970  -18.723    2.00  43116  13.0%  41546.847  12.5%  1569.153  7.698  22.970  7.651    910  1.00  19106  5.8%  17438.779  5.3%  1667.221  12.625  22.402  12.432    2.00  94079  28.4%  95746.221  28.9%  -1667.221  -5.388  -22.402  -5.404		56	1.00	1274	.4%	1459.536	.4%	-185.536	-4.856	-5.424	-4.965
78  1.00  5998  1.8%  7567.153  2.3%  -1569.153  -18.038  -22.970  -18.723    2.00  43116  13.0%  41546.847  12.5%  1569.153  7.698  22.970  7.651    910  1.00  19106  5.8%  17438.779  5.3%  1667.221  12.625  22.402  12.432    2.00  94079  28.4%  95746.221  28.9%  -1667.221  -5.388  -22.402  -5.404			2.00	8199	2.5%	8013.464	2.4%	185.536	2.073	5.424	2.065
2.00  43116  13.0%  41546.847  12.5%  1569.153  7.698  22.970  7.651    910  1.00  19106  5.8%  17438.779  5.3%  1667.221  12.625  22.402  12.432    2.00  94079  28.4%  95746.221  28.9%  -1667.221  -5.388  -22.402  -5.404		78	1.00	5998	1.8%	7567.153	2.3%	-1569.153	-18.038	-22.970	-18.723
910  1.00  19106  5.8%  17438.779  5.3%  1667.221  12.625  22.402  12.432    2.00  94079  28.4%  95746.221  28.9%  -1667.221  -5.388  -22.402  -5.404			2.00	43116	13.0%	41546.847	12.5%	1569.153	7.698	22.970	7.651
2.00 94079 28.4% 95746.221 28.9% -1667.221 -5.388 -22.402 -5.404		910	1.00	19106	5.8%	17438.779	5.3%	1667.221	12.625	22.402	12.432
			2.00	94079	28.4%	95746.221	28.9%	-1667.221	-5.388	-22.402	-5.404

Cell Counts and Residuals<sup>a,b</sup>

a. Model: Poisson

b. Design: Constant + movie\_n + rating + sex\_n + movie\_n \* rating + movie\_n \* sex\_n

Parameter Estimates <sup>b,c</sup>										
					95% Confidence Interval					
Parameter	Estimate	Std. Error	Z	Sig.	Lower Bound	Upper Bound				
Constant	11.469	.003	3656.315	.000	11.463	11.476				
$[movie_n = 1.00]$	230	.005	-48.520	.000	239	221				
$[movie_n = 2.00]$	$0^{a}$									
[rating = 12]	-2.818	.013	-225.062	.000	-2.843	-2.794				
[rating = 34]	-3.735	.019	-191.873	.000	-3.773	-3.697				
[rating = 56]	-2.481	.011	-231.923	.000	-2.502	-2.460				
[rating = 78]	835	.005	-154.512	.000	845	824				
[rating = 910]	$0^{a}$					•				

$[sex_n = 1.00]$	-1.703	.007	-261.736	.000	-1.716	-1.690
$[sex_n = 2.00]$	$0^{a}$					
[movie_n = 1.00] * [rating = 12]	952	.026	-36.357	.000	-1.004	901
[movie_n = 1.00] * [rating = 34]	129	.031	-4.165	.000	190	068
[movie_n = 1.00] * [rating = 56]	.180	.016	11.483	.000	.149	.211
[movie_n = 1.00] * [rating = 78]	.399	.008	51.683	.000	.384	.414
[movie_n = 1.00] * [rating = 910]	$0^{a}$					
[movie_n = 2.00] * [rating = 12]	0 <sup>a</sup>					
[movie_n = 2.00] * [rating = 34]	0 <sup>a</sup>					
[movie_n = 2.00] * [rating = 56]	0 <sup>a</sup>					
[movie_n = 2.00] * [rating = 78]	O <sup>a</sup>					
[movie_n = 2.00] * [rating = 910]	$0^{a}$					
$[movie_n = 1.00] * [sex_n = 1.00]$	576	.011	-52.282	.000	598	554
$[movie_n = 1.00] * [sex_n = 2.00]$	$0^{a}$					
$[movie_n = 2.00] * [sex_n = 1.00]$	0 <sup>a</sup>					
$[movie_n = 2.00] * [sex_n = 2.00]$	0 <sup>a</sup>					

a. This parameter is set to zero because it is redundant.

b. Model: Poisson

c. Design: Constant + movie\_n + rating + sex\_n + movie\_n \* rating + movie\_n \* sex\_n

# WE WISH TO TEST IF THE RATING IS INDEPENDENT FROM THE RATERS SEX, GIVEN DIFFERENT MOVIES. (THAT IS, WE WISH TO TEST IF IS IT TRUE THAT THE RELATIONSHIP BETWEEN RATING AND SEX VANISHES WHEN WE HOLD MOVIE CONSTANT.)

Ho:  $\lambda_{jk}^{\text{ratesex}} = \lambda_{ijk}^{\text{ratesexmovie}} = 0$  Ha: Ho not true,  $\alpha = 0.05$ 

Reject Ho. (Pearson  $\chi^2 = 1788.9$ , df =8, pvalue <.001, Likelihood Ration = 1599.9, df =8, p-value<0.001). Sex and rating are not independent, given the movie, and thus sex has a significant effect on the rating of a movie. (This follows from our previous result that the relationship between sex, rating, and movie is not homogenous)

# **<u>3 WAY TABLE: SUPERVISOR, WORKER SATISFACTION</u>**

Coding: manage\_n (0 - bad, 1 - good) Superv\_n(0 - low, 1-high) Worker\_n(0 - low, 1 - high)

File name: Satisfactionrecode.sav

WILL DO FOUR MODELS SATURATED MODEL (ALL MAIN, 2 WAY, AND 3 WAY EFFECTS) HOMOGENOUS MODEL (ALL MAIN AND 2 WAY EFFECTS) CONDITIONAL INDEPENDENCE MODEL (ALL MAIN EFFECT, MW INTERACTION AND MS INTERACTION) INDEPENDENT MODEL (MAIN EFFECTS ONLY)

# SATURATED MODEL

Data>Weight Cases Weight Cases by: freq Ok

Analysis>Loglinear>General Factor(s) Manage\_n Superv\_n Worker\_n Model: Custom Put manage\_n, superv\_n, and worker\_n in as main effects terms Put manage n\*superv n, manage n\*worker n, and superv n\*worker n in as 2way interaction terms Put manage\_n\*surperv\_n\*worker\_n in as a 3 way interaction term Continue **Options:** Display: Include Frequencies, Residuals, Design matrix and estimates Plot: Include Adjusted Residuals, Normal probability for adjusted Continue Save: Predicted Values Continue Ok

Below you will find the Goodness of Fit Test, Cell Counts and Residuals and Parameter Estimates tables obtained from running the commands above. Because this is the saturated model, SPSS adds .5 to all cell counts. This default is specified as Delta = .5 in the Criteria area of the options window.

Goodness-of-Fit Testsa,
-------------------------

	Value	df	Sig.
Likelihood Ratio	.000	0	•
Pearson Chi-Square	.000	0	
	MILDI	-	

a. Model: Poisson

b. Design: Constant + manage\_n + superv\_n + worker\_n + manage\_n \* superv\_n + manage\_n \* worker\_n + superv\_n \* worker\_n + manage\_n \* superv\_n \* worker\_n

Cell Counts and Residualsa,b

			Obse	Observed		Expected		Standardized	Adjusted	
manage_n	superv_n	worker_n	Count	%	Count	%	Residual	Residual	Residual	Deviance
.00	.00	.00	103.500	14.4%	103.500	14.4%	.000	.000	•	.000
		1.00	87.500	12.2%	87.500	12.2%	.000	.000	•	.000
	1.00	.00	32.500	4.5%	32.500	4.5%	.000	.000	.000	.000
		1.00	42.500	5.9%	42.500	5.9%	.000	.000	.000	.000
1.00	.00	.00	59.500	8.3%	59.500	8.3%	.000	.000	.000	.000
		1.00	109.500	15.2%	109.500	15.2%	.000	.000	.000	.000
	1.00	.00	78.500	10.9%	78.500	10.9%	.000	.000	•	.000
	-	1.00	205.500	28.6%	205.500	28.6%	.000	.000	.000	.000

a. Model: Poisson

 $b. Design: Constant + manage_n + superv_n + worker_n + manage_n * superv_n + manage_n * worker_n + superv_n * worker_n + manage_n * superv_n * superv_n * superv_n * superv_n * superv_n * superv_n$ 

		Paramete	r Estimatesb,c			
					95% Confide	ence Interval
Parameter	Estimate	Std. Error	Z	Sig.	Lower Bound	Upper Bound
Constant	5.325	.070	76.342	.000	5.189	5.462
$[manage_n = .00]$	-1.576	.169	-9.352	.000	-1.906	-1.246
$[manage_n = 1.00]$	0a		•			•
$[superv_n = .00]$	630	.118	-5.321	.000	861	398
[superv_n = 1.00]	0a		•		•	•
[worker_n = .00]	962	.133	-7.253	.000	-1.222	702
[worker_n = 1.00]	0a		•		•	•
[manage_n = .00] * [superv_n =	1.352	.221	6.109	.000	.918	1.785
.00]						
[manage_n = .00] * [superv_n =	0a		•			•
1.00]						
[manage_n = 1.00] * [superv_n =	0a		•			•
.00]						
[manage_n = 1.00] * [superv_n =	0a		•		•	•
1.00]						
[manage_n = .00] * [worker_n =	.694	.268	2.588	.010	.169	1.220
.00]						
[manage_n = .00] * [worker_n =	0a		•		•	•
1.00]						
[manage_n = 1.00] * [worker_n =	0a		•		•	•
.00]						
[manage_n = 1.00] * [worker_n =	0a	•	•		•	•
1.00]						
[superv_n = .00] * [worker_n =	.352	.209	1.689	.091	057	.761
.00]						
[superv_n = .00] * [worker_n =	0a	•	•	•	•	•
1.00]						
[superv_n = 1.00] * [worker_n =	0a	•	•	•	•	•
.00]						
[superv_n = 1.00] * [worker_n =	0a	•	•	•	•	•
1.00]						
[manage_n = .00] * [superv_n =	.084	.345	.243	.808	592	.760
.00] * [worker_n = .00]						
[manage_n = .00] * [superv_n =	0a	•	•	•	•	•
$.00] * [worker_n = 1.00]$	0					
$[\text{manage}_n = .00] * [\text{superv}_n = .00]$	Ua	•	•	•	•	•
$1.00] * [worker_n = .00]$	0					
$[manage_n = .00] * [superv_n = 1.00] * [$	Ua	•	•	•	•	•
$1.00 ] * [worker_n = 1.00]$	0-					
$[manage_n = 1.00] * [superv_n = 0.0] * [superv_n = 0.0]$	Ua	•	•	•	•	•
$.00] = [worker_n = .00]$	00					
$[\text{Inamage_II} - 1.00] * [\text{superv_II} = 0.0] * [\text{worker } n = 1.00]$	va	•	•	•	•	•
[monogo n = 1.00] * [superv n = 1.00]	0.0					
$[\text{Inamage_II} - 1.00] * [\text{superv_II} = 1.00] * [\text{worker} n = 0.01]$	va	•	•	•	•	•
$1.00 ] \cdot [w01 \text{Ker}_{II} = .00]$	0.5					
1.00] * [worker n - 1.00]	va	•	•		•	•
1.00] [WOLKEL_H = 1.00]		•••••	h			

 $\label{eq:c.Design: Constant + manage_n + superv_n + worker_n + manage_n * superv_n + manage_n * worker_n + superv_n * worker_n + manage_n * superv_n *$ 

If you reset Delta = 0, SPSS will return a Cell counts and Residuals table without the extra .5 added. The Goodness of Fit test results are not affected, but the Parameter Estimates will change slightly.

# goodness-of-Fit Tests<sup>a,b</sup>

	Value	df	Sig.
Likelihood Ratio	.000	0	
Pearson Chi-	.000	0	
Square			

a. Model: Poisson

b. Design: Constant + manage\_n + superv\_n +

 $worker\_n + manage\_n * superv\_n + manage\_n *$ 

worker\_n + superv\_n \* worker\_n + manage\_n \*

superv\_n \* worker\_n

	-	-	Obse	erved	Expe	ected		Standardi		
manage	superv	worker					Residu	zed	Adjusted	Devian
_n	_n	_n	Count	%	Count	%	al	Residual	Residual	ce
.00	.00	.00	103	14.4%	103.00	14.4%	.000	.000		.000
					0					
		1.00	87	12.2%	87.000	12.2%	.000	.000	.000	.000
	1.00	.00	32	4.5%	32.000	4.5%	.000	.000	.000	.000
		1.00	42	5.9%	42.000	5.9%	.000	.000	.000	.000
1.00	.00	.00	59	8.3%	59.000	8.3%	.000	.000	.000	.000
		1.00	109	15.2%	109.00	15.2%	.000	.000	.000	.000
					0					
	1.00	.00	78	10.9%	78.000	10.9%	.000	.000	.000	.000
		1.00	205	28.7%	205.00	28.7%	.000	.000	.000	.000
					0					

#### Cell Counts and Residuals<sup>a,b</sup>

a. Model: Poisson

b. Design: Constant + manage\_n + superv\_n + worker\_n + manage\_n \* superv\_n + manage\_n \* worker\_n + superv\_n \* worker\_n + manage\_n \* superv\_n \* worker\_n

				-	95% Confidence Interval	
Parameter	Estimate	Std. Error	Z	Sig.	Lower Bound	Upper Bound
Constant	5.323	.070	76.214	.000	5.186	5.460
$[manage_n = .00]$	-1.585	.169	-9.360	.000	-1.917	-1.253
$[manage_n = 1.00]$	0 <sup>a</sup>					
$[superv_n = .00]$	632	.119	-5.329	.000	864	399
$[superv_n = 1.00]$	0 <sup>a</sup>					
[worker_ $n = .00$ ]	966	.133	-7.263	.000	-1.227	706
$[worker_n = 1.00]$	$0^{a}$					
[manage_n = .00] * [superv_n =	1.360	.222	6.121	.000	.924	1.795
.00]						
[manage_n = .00] * [superv_n =	0 <sup>a</sup>					
1.00]						
$[manage_n = 1.00] * [superv_n =$	0 <sup>a</sup>					
.00]						
[manage_n = 1.00] * [superv_n =	$0^{a}$					
1.00]						
$[manage_n = .00] * [worker_n =$	.694	.270	2.574	.010	.166	1.223
.00]						
$[manage_n = .00] * [worker_n =$	0 <sup>a</sup>		-			
1.00]						
$[manage_n = 1.00] * [worker_n =$	0 <sup>a</sup>					
.00]						
$[manage_n = 1.00] * [worker_n =$	0 <sup>a</sup>		•			
1.00]						
$[superv_n = .00] * [worker_n = .00]$	.352	.209	1.684	.092	058	.763
$[superv_n = .00] * [worker_n =$	0 <sup>a</sup>		•		•	•
1.00]						
$[superv_n = 1.00] * [worker_n =$	0 <sup>a</sup>		•			
.00]						
$[superv_n = 1.00] * [worker_n =$	0 <sup>a</sup>			•		
1.00]						
$[manage_n = .00] * [superv_n =$	.088	.347	.255	.799	591	.767
$.00] * [worker_n = .00]$						
$[manage_n = .00] * [superv_n =$	0 <sup>a</sup>			•		
$.00] * [worker_n = 1.00]$						
$[manage_n = .00] * [superv_n =$	0ª	•	•	•	•	•
$1.00$ ] * [worker_n = .00]						
$[manage_n = .00] * [superv_n = .00]$	0ª	•	•	•	•	•
$[1.00] * [worker_n = 1.00]$	03					
$[manage_n = 1.00] * [superv_n = 0.01]$	0"		•		•	
$.00] * [worker_n = .00]$	03					
$[manage_n = 1.00] * [superv_n = 1.00]$	0"		•		•	•
$.00] * [worker_n = 1.00]$	<u></u>					
$[manage_n = 1.00] * [superv_n = 1.00] * [superv_n = 0.01]$	$0^{a}$	•	•	•	•	
$1.00] + [worker_n = .00]$	0.0					
$[manage_n = 1.00] * [superv_n = 1.00]$	$0^{a}$		•	•	•	
$1.00] \ (worker_n = 1.00]$						

Parameter Estimates<sup>b,c</sup>

a. This parameter is set to zero because it is redundant.

b. Model: Poisson

c. Design: Constant + manage\_n + superv\_n + worker\_n + manage\_n \* superv\_n + manage\_n \* worker\_n + superv\_n \* worker\_n + manage\_n \* superv\_n \* worker\_n

# HOMOGENOUS MODEL (INCLUDES MAIN EFFECTS AND ALL 2 WAY EFFECTS)

Data>Weight Cases Weight cases by Frequency Variable: freq

# OK

Analysis>Loglinear>General Factor(s) Manage\_n Superv\_n Worker\_n Model: Custom Put manage\_n, superv\_n, and worker\_n in as main effects terms Put manage\_n\*superv\_n, manage\_n\*worker\_n, and superv\_n\*worker\_n in as 2way interaction terms Continue **Options:** Display: Include Frequencies, Residuals, Design matrix and estimates Plot: Include Adjusted Residuals, Normal probability for adjusted Continue Save: Predicted Values Continue OK

## Goodness-of-Fit Tests<sup>a,b</sup>

	Value	df	Sig.
Likelihood Ratio	.065	1	.799
Pearson Chi-	.065	1	.799
Square			

a. Model: Poisson

b. Design: Constant + manage\_n + superv\_n + worker\_n + manage\_n \* superv\_n + manage\_n \* worker\_n + superv\_n \* worker\_n

	-		Obse	rved	Expe	cted		Standardized	Adjusted			
manage_n	superv_n	worker_n	Count	%	Count	%	Residual	Residual	Residual	Deviance		
.00	.00	.00	103	14.4%	102.264	14.3%	.736	.073	.255	.073		
		1.00	87	12.2%	87.736	12.3%	736	079	255	079		
	1.00	.00	32	4.5%	32.736	4.6%	736	129	255	129		
		1.00	42	5.9%	41.264	5.8%	.736	.115	.255	.114		
1.00	.00	.00	59	8.3%	59.736	8.4%	736	095	255	095		
		1.00	109	15.2%	108.264	15.1%	.736	.071	.255	.071		
	1.00	.00	78	10.9%	77.264	10.8%	.736	.084	.255	.084		
		1.00	205	28.7%	205.736	28.8%	736	051	255	051		

Cell Counts and Residuals<sup>a,b</sup>

a. Model: Poisson

-		-	Obse	erved	Expe	ected		Standardized	Adjusted	
manage_n	superv_n	worker_n	Count	%	Count	%	Residual	Residual	Residual	Deviance
.00	.00	.00	103	14.4%	102.264	14.3%	.736	.073	.255	.073
		1.00	87	12.2%	87.736	12.3%	736	079	255	079
	1.00	.00	32	4.5%	32.736	4.6%	736	129	255	129
		1.00	42	5.9%	41.264	5.8%	.736	.115	.255	.114
1.00	.00	.00	59	8.3%	59.736	8.4%	736	095	255	095
		1.00	109	15.2%	108.264	15.1%	.736	.071	.255	.071
	1.00	.00	78	10.9%	77.264	10.8%	.736	.084	.255	.084
		1.00	205	28.7%	205.736	28.8%	736	051	255	051

Cell Counts and Residuals<sup>a,b</sup>

 $b. Design: Constant + manage_n + superv_n + worker_n + manage_n * superv_n + manage_n * worker_n + superv_n * worker_n + manage_n * worker_n + manage_n$ 

			-		95% Confide	ence Interval
Parameter	Estimate	Std. Error	Z	Sig.	Lower Bound	Upper Bound
Constant	5.327	.068	78.002	.000	5.193	5.460
$[manage_n = .00]$	-1.607	.148	-10.826	.000	-1.897	-1.316
$[manage_n = 1.00]$	0 <sup>a</sup>					
$[superv_n = .00]$	642	.112	-5.757	.000	861	423
$[superv_n = 1.00]$	0 <sup>a</sup>					
$[worker_n = .00]$	979	.123	-7.955	.000	-1.221	738
$[worker_n = 1.00]$	0 <sup>a</sup>					
[manage_n = .00] * [superv_n =	1.396	.171	8.189	.000	1.062	1.731
.00]						
[manage_n = .00] * [superv_n = 1.00]	$0^{a}$	•				
$[manage_n = 1.00] * [superv_n =$	0 <sup>a</sup>					
.00]	0.0					
$[manage_n = 1.00] * [superv_n = 1.00]$	0ª					
[manage_n = .00] * [worker_n = .00]	.748	.169	4.423	.000	.416	1.079
[manage_n = .00] * [worker_n = 1.00]	$0^{a}$					
[manage_n = 1.00] * [worker_n = .00]	$0^{a}$					
[manage_n = 1.00] * [worker_n = 1.00]	$0^{a}$					
$[superv_n = .00] * [worker_n = .00]$	.385	.167	2.309	.021	.058	.711
[superv_n = .00] * [worker_n =	0 <sup>a</sup>					
1.00]						
$[superv_n = 1.00] * [worker_n =$	0 <sup>a</sup>					
.00]						
[superv_n = 1.00] * [worker_n =	$0^{a}$		•			•
1.00]						

#### Parameter Estimates<sup>b,c</sup>

a. This parameter is set to zero because it is redundant.

b. Model: Poisson

 $c. \ Design: Constant + manage_n + superv_n + worker_n + manage_n * superv_n + manage_n * worker_n + superv_n * worker_n + manage_n * worker_n + worker_n + worker_n + worker_n + worker_n + worker_$ 

# CONDITIONAL INDEPENDENCE MODEL (CONDITIONED ON SUPERV-WORKER) (INCLUDES MAIN EFFECTS, MW AND MS TWO WAY INTERACTION TERMS)

Weight cases by Frequency Variable: freq OK

Analysis>Loglinear>General Factor(s) Manage\_n Superv\_n Worker\_n Model: Custom Put manage\_n, superv\_n, and worker\_n in as main effects terms Put manage\_n\*superv\_n, and manage\_n\*worker\_n, in as 2way interaction terms Continue **Options:** Display: Include Frequencies, Residuals, Design matrix and estimates Plot: Include Adjusted Residuals, Normal probability for adjusted Continue Save: Predicted Values Continue Ok

Goodness-of-Fit Tests<sup>a,b</sup>

	Value	df	Sig.
Likelihood Ratio	5.387	2	.068
Pearson Chi-	5.410	2	.067
Square			

a. Model: Poisson

b. Design: Constant + manage\_n + superv\_n + worker\_n + manage\_n \* superv\_n + manage\_n \* worker\_n

			Obse	Observed		Expected		Standardized	Adjusted	
manage_n	superv_n	worker_n	Count	%	Count	%	Residual	Residual	Residual	Deviance
.00	.00	.00	103	14.4%	97.159	13.6%	5.841	.593	1.601	.587
		1.00	87	12.2%	92.841	13.0%	-5.841	606	-1.601	613
	1.00	.00	32	4.5%	37.841	5.3%	-5.841	950	-1.600	976
		1.00	42	5.9%	36.159	5.1%	5.841	.971	1.600	.947
1.00	.00	.00	59	8.3%	51.033	7.1%	7.967	1.115	1.687	1.088
		1.00	109	15.2%	116.967	16.4%	-7.967	737	-1.687	745

Cell Counts and Residuals<sup>a,b</sup>

1.00 .00	78	10.9%	85.967	12.0%	-7.967	859	-1.687	873
1.00	205	28.7%	197.033	27.6%	7.967	.568	1.684	.564

 $b. \ Design: Constant + manage_n + superv_n + worker_n + manage_n * superv_n + manage_n * worker_n$ 

Parameter Estimates <sup>o,e</sup>												
					95% Confide	ence Interval						
Parameter	Estimate	Std. Error	Z	Sig.	Lower Bound	Upper Bound						
Constant	5.283	.067	78.772	.000	5.152	5.415						
$[manage_n = .00]$	-1.695	.148	-11.440	.000	-1.986	-1.405						
$[manage_n = 1.00]$	$0^{a}$				•							
$[superv_n = .00]$	521	.097	-5.355	.000	712	331						
$[superv_n = 1.00]$	$0^{a}$				•							
$[worker_n = .00]$	829	.102	-8.101	.000	-1.030	629						
$[worker_n = 1.00]$	0 <sup>a</sup>											
[manage_n = .00] * [superv_n =	1.464	.168	8.713	.000	1.135	1.794						
.00]												
[manage_n = .00] * [superv_n =	0 <sup>a</sup>											
1.00]												
[manage_n = 1.00] * [superv_n =	0 <sup>a</sup>											
.00]												
$[manage_n = 1.00] * [superv_n =$	0 <sup>a</sup>											
1.00]												
[manage_n = .00] * [worker_n =	.875	.160	5.464	.000	.561	1.189						
.00]												
$[manage_n = .00] * [worker_n =$	0 <sup>a</sup>											
1.00]												
[manage_n = 1.00] * [worker_n =	0 <sup>a</sup>											
.00]												
[manage_n = 1.00] * [worker_n =	O <sup>a</sup>											
1.00]												

a. This parameter is set to zero because it is redundant.

b. Model: Poisson

c. Design: Constant + manage\_n + superv\_n + worker\_n + manage\_n \* superv\_n + manage\_n \* worker\_n

# **INDEPENDENCE MODEL (ALL MAIN EFFECTS)**

Data>Weight Cases Weight cases by Frequency Variable: freq OK

Analysis>Loglinear>General Factor(s) Manage\_n Superv\_n Worker\_n Model: Custom Put manage\_n, superv\_n, and worker\_n in as main effects terms Continue Options: Display: Include Frequencies, Residuals, Design matrix and estimates Plot: Include Adjusted Residuals, Normal probability for adjusted Continue Save: Predicted Values Continue Ok

Goodness-of-Fit Tests<sup>a,b</sup>

	Value	df	Sig.
Likelihood Ratio	117.997	4	.000
Pearson Chi-	128.086	4	.000
Square			

a. Model: Poisson

b. Design: Constant + manage\_n + superv\_n +

worker\_n

	-	-	Obse	erved	Expe	ected		Standardi		
manage	superv	worker					Residu	zed	Adjusted	Devian
_n	_n	_n	Count	%	Count	%	al	Residual	Residual	ce
.00	.00	.00	103	14.4%	50.286	7.0%	52.714	7.434	9.404	6.502
		1.00	87	12.2%	81.899	11.5%	5.101	.564	.793	.558
	1.00	.00	32	4.5%	50.145	7.0%	-	-2.562	-3.240	-2.746
							18.145			
		1.00	42	5.9%	81.670	11.4%	-	-4.390	-6.171	-4.845
							39.670			
1.00	.00	.00	59	8.3%	85.905	12.0%	-	-2.903	-4.130	-3.078
							26.905			
		1.00	109	15.2%	139.91	19.6%	-	-2.613	-4.270	-2.720
					1		30.911			
	1.00	.00	78	10.9%	85.665	12.0%	-7.665	828	-1.177	841
		1.00	205	28.7%	139.52	19.5%	65.480	5.544	9.051	5.178
					0					

Cell Counts and Residuals<sup>a,b</sup>

a. Model: Poisson

	-	-	Obse	erved	Expe	ected		Standardi		
manage	superv	worker					Residu	zed	Adjusted	Devian
_n	_n	_n	Count	%	Count	%	al	Residual	Residual	ce
.00	.00	.00	103	14.4%	50.286	7.0%	52.714	7.434	9.404	6.502
		1.00	87	12.2%	81.899	11.5%	5.101	.564	.793	.558
	1.00	.00	32	4.5%	50.145	7.0%	-	-2.562	-3.240	-2.746
							18.145			
		1.00	42	5.9%	81.670	11.4%	-	-4.390	-6.171	-4.845
							39.670			
1.00	.00	.00	59	8.3%	85.905	12.0%	-	-2.903	-4.130	-3.078
							26.905			
		1.00	109	15.2%	139.91	19.6%	-	-2.613	-4.270	-2.720
					1		30.911			
	1.00	.00	78	10.9%	85.665	12.0%	-7.665	828	-1.177	841
		1.00	205	28.7%	139.52	19.5%	65.480	5.544	9.051	5.178
					0					

Cell Counts and Residuals<sup>a,b</sup>

b. Design: Constant + manage\_n + superv\_n + worker\_n

					95% Confide	ence Interval
		Std.			Lower	Upper
Parameter	Estimate	Error	Z	Sig.	Bound	Bound
Constant	4.938	.067	73.791	.000	4.807	5.069
$[manage_n = .00]$	536	.077	-6.911	.000	687	384
[manage_n =	0 <sup>a</sup>	•	•			
1.00]						
$[superv_n = .00]$	.003	.075	.037	.970	144	.149
$[superv_n = 1.00]$	0 <sup>a</sup>					
$[worker_n = .00]$	488	.077	-6.332	.000	639	337
[worker_n =	0 <sup>a</sup>	•	•			
1.00]						

# Parameter Estimates<sup>b,c</sup>

a. This parameter is set to zero because it is redundant.

b. Model: Poisson

ſ

c. Design: Constant + manage\_n + superv\_n + worker\_n

We can use the data above to create a table of actual and expected cell counts for our models of interest. These expected cell counts are useful in obtaining estimated odds ratios, as will be illustrated below.

Expected			Expected
----------	--	--	----------

Page

Manage	Superv	Worker	Observed	Saturated	Homogenous	Condl	Independ
						Indep	
						Superv-	
						Wrkr	
Bad(0)	Low(0)	Low(0)	103.00	103.00	102.26	97.16	50.29
Bad(0)	Low(0)	High(1)	87.00	87.00	87.74	92.84	81.90
Bad(0)	High(1)	Low(0)	32.00	32.00	32.74	37.84	50.15
Bad(0)	High(1)	High(1)	42.00	42.00	41.26	36.16	81.67
High(1)	Low(0)	High(1)	59.00	59.00	597.74	51.03	85.91
High(1)	Low(0)	Low(0)	109.00	109.00	108.26	116.9	139.91
High(1)	High(1)	High(1)	78.00	78.00	77.26	85.97	85.67
High(1)	High(1)	Low(0)	205.00	205.00	205.74	197.03	139.52

# NOTE; THESE ARE CORRECT, AS THEY WORK IN CALCULATING ODDS.

We can check for the fitness of these models by creating a table that takes information from the Goodness of Fit tables above.

Model	Like Rat $\chi^2$	Df	P-value	Pearson $\chi^2$	Df	P-value
Saturated	0	0	1	0	0	1
Homogenous	0.065	1	.799	0.065	1	.799
Cond Ind	5.387	2	0.068	5.410	2	0.067
(Superv-						
Wrkr)						
Independence	117.997	4	< 0.001	128.086	4	.001

We see that the conditional independent model is acceptable (with a p-value of 0.068 it is not rejected at the 5% significance level). This model has less parameters than the homogenous model, which makes it desirable.

We bring down the cell count and residual table from above and take another look at it for this conditional independent model. It does not indicate that there are any high standardized residuals to concern us.

	-	-	Obse	erved	Expe	ected		Standardized	Adjusted	
manage_n	superv_n	worker_n	Count	%	Count	%	Residual	Residual	Residual	Deviance
.00	.00	.00	59	8.3%	51.033	7.1%	7.967	1.115	1.687	1.088
		1.00	109	15.2%	116.967	16.4%	-7.967	737	-1.687	745
	1.00	.00	78	10.9%	85.967	12.0%	-7.967	859	-1.687	873
		1.00	205	28.7%	197.033	27.6%	7.967	.568	1.684	.564
1.00	.00	.00	103	14.4%	97.159	13.6%	5.841	.593	1.601	.587
		1.00	87	12.2%	92.841	13.0%	-5.841	606	-1.601	613
	1.00	.00	32	4.5%	37.841	5.3%	-5.841	950	-1.600	976
		1.00	42	5.9%	36.159	5.1%	5.841	.971	1.600	.947

Cell Counts and Residuals<sup>a,b</sup>

a. Model: Poisson

 $b. \ Design: Constant + manage_n + superv_n + worker_n + manage_n * superv_n + manage_n * worker_n + manage_n + superv_n + superv_n + manage_n + superv_n + superv_$ 

We can also test Ho:  $\lambda_{ij}^{SW} = 0$  versus Ha:  $\lambda_{ij}^{SW} \neq 0$ ,  $\alpha = 0.05$ Test statistic =  $-2(L_o - L_1) = 5.387 - 0.065$  (both from output) = 5.321, df = 2-1 = 1From Table 7 in the textbook, 0.01 < p-value < 0.025We reject Ho. The data provides evidence at the 5% significance level that this interaction term,  $\lambda_{ij}^{SW} \neq 0$ . It appears from here that  $\lambda_{ij}^{SW}$  should remain in the model.

We have contradictory results. We will have to decide whether we might prefer the parsimonious model given that the residuals do not indicate trouble.

We can also check AIC = -2(Log likelihood - #parameters) for both models, again using -2Log Likelihoods from the output. Below we can see that it is smaller for the homogenous model, thus favouring the homogenous model.

Model	AIC
Homogenous	.065 - 2(7) = -13.94
Cond Indep (Super	5.387 - 2(6) = -6.61
Worker)	

# **Fitted Odds Ratios:**

Before we use the data above to look at how to obtain fitted odd ratios, we will take a bit of a background theoretical refresher to help us out.

Loglinear models use  $\mu_{ij} = n\pi_{ij}$  rather than  $\pi_{ij}$ .

With 2 variables X and Y, and two levels to each variable, (X=0,1 and Y=0,1), we can write

Log  $\mu_{ij} = \lambda + \lambda_i^X + \lambda_j^Y$ , i =0,1 and j = 0,1 For our work, we have assumed that the count is a Poisson variable. In this case, when we look at the odds we note:

 $\Theta_{XY} = \frac{n_{00}n_{11}}{n_{01}n_{10}} = \frac{\pi_{00}\pi_{11}}{\pi_{01}\pi_{10}} = \frac{\mu_{00}\mu_{11}}{\mu_{01}\mu_{10}}$  and our table of interest might be set up this way

	Y(0)	Y(1)
X(0)	n <sub>00</sub>	n <sub>01</sub>
X(1)	<b>n</b> <sub>10</sub>	<b>n</b> <sub>11</sub>

We can't go wrong with the odds if we say:

The odds of observing Y level 0 are #times higher for X level 0 than X level 1.

Now let us look at a case with 3 variables, X, Y, and Z, each with two levels (0 and 1).

 $Log \; \mu_{ijk} = \lambda + \lambda_i^X + \lambda_j^Y + \lambda_k^Z + \lambda_{ij}^{XY} + \lambda_{ik}^{XZ} + \lambda_{jk}^{YZ} + \lambda_{ijk}^{XYZ} \quad i = 0, 1, \; j = 0, 1, \; and \; k = 0, 1 \\ (1 + \lambda_i^X) + \lambda_i^X +$ 

If we look at levels 0 and 1 separately for Z, we can have two tables of interest, viz: Z(0)

	Y(0)	Y(1)						
X(0)	n <sub>00</sub>	n <sub>01</sub>						
X(1)	<b>n</b> <sub>10</sub>	<b>n</b> <sub>11</sub>						
Here w	e look at	t Log(µi	$_{ij0} = \lambda + \lambda$	$L_i^X + \lambda_j^Y + \lambda_0$	$^{Z} + \lambda_{ij}^{XY} + \lambda_{ij}^{XY}$	$\lambda_{i0}^{XZ} + \lambda_{j0}^{YZ} +$	$-\lambda_{ij0}^{XYZ}$	i=0,1 and $j=0$
$\Theta_{\mathbf{V}\mathbf{V}(0)}$	$=\frac{n_{00(0)}n_1}{n_1}$	$\frac{\pi_{11(0)}}{\pi_{11(0)}} = \frac{\pi_{11(0)}}{\pi_{11(0)}}$	$\pi_{11(0)}$	$\mu_{00(0)}\mu_{11(0)}$				
$O_{X1}(0)$	$n_{01(0)}n_1$	$\pi_{0}(0) = \pi_{0}(0)$	$\pi_{100}\pi_{10(0)}$	$\mu_{01(0)}\mu_{10(0)}$				

We would say that the odds of observing Y level 0 are #times higher for X level 0 than X level 1 (conditioned on Z=0).

. ,
-----

	Y(0)	Y(1)
X(0)	<b>n</b> 00	<b>n</b> 01
X(1)	<b>n</b> <sub>10</sub>	<b>n</b> <sub>11</sub>

Here we look at  $\text{Log}(\mu_{ij1}) = \lambda + \lambda_i^X + \lambda_j^Y + \lambda_1^Z + \lambda_{ij}^{XY} + \lambda_{i1}^{XZ} + \lambda_{i1}^{YZ} + \lambda_{ij1}^{XYZ}$  i= 0,1 and j=0,1 With corresponding odds ratios of interest expressed as:

 $\Theta_{XY(1)} = \frac{n_{00(1)}n_{11(1)}}{n_{01(1)}n_{10(1)}} = \frac{\pi_{00(1)}\pi_{11(1)}}{\pi_{01(1)}\pi_{10(1)}} = \frac{\mu_{00(1)}\mu_{11(1)}}{\mu_{01(1)}\mu_{10(1)}}$ 

We would say that the odds of observing Y level 0 are #times higher for X level 0 than X level 1 (conditioned on Z=1).

We can create two 2x2 tables in Y and Z if we look at levels 0 and 1 separately for X. We can create two 2x2 tables in X and Z if we look at levels 0 and 1 separately for Y.

For our example, we are considering various loglinear Poisson models, as follows.

# **Saturated**

# Homogenous

# **Conditional Independent (Supervisor-Worker)**

# Independent

Equations for these are as follows:

Saturated : Log  $\mu_{ijk} = \lambda + \lambda_i^M + \lambda_j^S + \lambda_k^W + \lambda_{ij}^{MS} + \lambda_{ik}^{MW} + \lambda_{jk}^{SW} + \lambda_{ijk}^{MSW}$ i=0,1, j=0,1, and k=0,1Homogen : Log  $\mu_{ijk} = \lambda + \lambda_i^M + \lambda_j^S + \lambda_k^W + \lambda_{ij}^{MS} + \lambda_{ik}^{MW} + \lambda_{jk}^{SW}$ i=0,1, j=0,1, and k=0,1CondI(SW): Log  $\mu_{ijk} = \lambda + \lambda_i^M + \lambda_j^S + \lambda_k^W + \lambda_{ij}^{MS} + \lambda_{ik}^{MW}$ i=0,1, i=0,1, and k=0,1Independ : Log  $\mu_{ijk} = \lambda + \lambda_i^M + \lambda_i^S + \lambda_k^W$ i=0,1, i=0,1, and k=0,1

Note that each of these e	equations has eight	possible values,	corresponding to
---------------------------	---------------------	------------------	------------------

Manage	Superv	Worker
Bad(0)	Low(0)	Low(0)
Bad(0)	Low(0)	High(1)
Bad(0)	High(1)	Low(0)
Bad(0)	High(1)	High(1)
High(1)	Low(0)	Low(0)

High(1)	Low(0)	High(1)
High(1)	High(1)	Low(0)
High(1)	High(1)	High(1)

Now let us consider some conditional odds that might be of interest for our models.

Model	M and S (condition on W)	M and W (condition on S)	S and W (condition on M)
Saturated – level 0	$\Theta_{MS(0)}$	$\Theta_{M(0)W}$	$\Theta_{(0)SW}$
Saturated – level 1	$\Theta_{MS(1)}$	$\Theta_{M(1)W}$	$\Theta_{(1)SW}$
Homogenous	$\Theta_{MS(0)} = \Theta_{MS(1)}$	$\Theta_{M(0)W} = \Theta_{M(1)W}$	$\Theta_{(0)SW} = \Theta_{(1)SW}$
Cond. Indep (Super-	$\Theta_{MS(0)} = \Theta_{MS(1)}$	$\Theta_{M(0)W} = \Theta_{M(1)W}$	$\Theta_{(0)SW} = \Theta_{(1)SW}$
Worker)			
Independence	$\Theta_{MS(0)} = \Theta_{MS(1)}$	$\Theta_{M(0)W} = \Theta_{M(1)W}$	$\Theta_{(0)SW} = \Theta_{(1)SW}$

For the Homogenous and the Conditional Independence Model, the conditional odds are the same for both levels of the variable on which we condition. We prove this for the Conditional Independence Model (Super-Worker) situation above for the M and S (condition on W) situation. Other proofs are analogous.

$$\begin{split} & \text{Log } \Theta_{\text{MS}(k)} = \text{Log } (\frac{\mu_{00}(k)\mu_{11}(k)}{\mu_{01}(k)\mu_{10}(k)}) = \text{log}\mu_{00k} + \text{log}\mu_{11k} - \text{log}\mu_{01k} - \text{log}\mu_{10k} = \\ & \lambda + \lambda_0^M + \lambda_0^S + \lambda_k^W + \lambda_{00}^{\text{MS}} + \lambda_{0k}^{\text{MW}} \\ & + \lambda + \lambda_1^M + \lambda_1^S + \lambda_k^W + \lambda_{11}^{\text{MS}} + \lambda_{1k}^{\text{MW}} \\ & -(\lambda + \lambda_0^M + \lambda_1^S + \lambda_k^W + \lambda_{01}^{\text{MS}} + \lambda_{0k}^{\text{MW}}) \\ & -(\lambda + \lambda_1^M + \lambda_0^S + \lambda_k^W + \lambda_{10}^{\text{MS}} + \lambda_{1k}^{\text{MW}}) \\ & = \lambda_{00}^{\text{MS}} + \lambda_{11}^{\text{MS}} - \lambda_{01}^{\text{MS}} - \lambda_{10}^{\text{MS}} \end{split}$$
which does not depend on k. So Log  $\Theta_{\text{MS}(0)} = \text{Log } \Theta_{\text{MS}(1)}$ , and exponentializing gives us  $\Theta_{\text{MS}(0)} = \Theta_{\text{MS}(1)}$ .

Furthermore, for us, this is equal to  $\lambda_{00}^{MS}$  because SPSS sets up the equation so that the only non-zero parameter estimates returned are when levels of the all categories involved in the estimated parameter are 0.

Further algebra will provide us with the following equations.

Model	M and S (condition on W)	M and W (condition on S)	S and W (condition on M)
Saturated – level 0	$\Theta_{\mathrm{MS}(0)} = e^{\lambda_{00}^{MS} + \lambda_{000}^{MSW}}$	$\Theta_{\mathrm{M}(0)\mathrm{W}} = e^{\lambda_{00}^{MW} + \lambda_{000}^{MSW}}$	$\Theta_{(0)SW} = e^{\lambda_{00}^{SW} + \lambda_{000}^{MSW}}$
Saturated – level 1	$\Theta_{\rm MS(1)} = e^{\lambda_{00}^{\rm MS}}$	$\Theta_{\mathrm{M}(1)\mathrm{W}} = e^{\lambda_{00}^{MW}}$	$\Theta_{(1)SW} = e^{\lambda_{00}^{SW}}$
Homogenous	$\Theta_{\rm MS(0)} = \Theta_{\rm MS(1)} = e^{\lambda_{00}^{MS}}$	$\Theta_{\mathrm{M}(0)\mathrm{W}} = \Theta_{\mathrm{M}(1)\mathrm{W}} = e^{\lambda_{00}^{MW}}$	$\Theta_{(0)SW} = \Theta_{(1)SW} = e^{\lambda_{00}^{SW}}$
Cond. Indep (Super-	$\Theta_{\rm MS(0)} = \Theta_{\rm MS(1)} = e^{\lambda_{00}^{MS}}$	$\Theta_{\mathrm{M}(0)\mathrm{W}} = \Theta_{\mathrm{M}(1)\mathrm{W}} = e^{\lambda_{00}^{MW}}$	1
Worker)			
Independence	1	1	1

For our models, our output provides the estimated parameters for the 0 levels as below.

	λ	$\lambda_0^{M}$	$\lambda_0^{S}$	$\lambda_0^{W}$	$\lambda_{00}^{MS}$	$\lambda_{00}^{MW}$	$\lambda_{00}^{SW}$	$\lambda_{000}^{MSW}$	
Saturated	5.323	-1.584	632	966	1.360	.694	.352	.088	
	λ	$\lambda_0^{M}$	$\lambda_0^{S}$	$\lambda_0^{W}$	$\lambda_{00}^{MS}$	$\lambda_{00}^{MW}$	$\lambda_{00}^{SW}$		S <sup>e</sup>
Homogen	5.327	-1.607	642	979	1.396	.748	.385		Pag

	λ	$\lambda_0^M$	$\lambda_0^{S}$	$\lambda_0^{W}$	$\lambda_{00}^{MS}$	$\lambda_{00}^{MW}$	
CondI(SW)	5.283	-1.695	521	829	1.464	.875	
	λ	$\lambda_0^{M}$	$\lambda_0^{S}$	$\lambda_0^{W}$			
Indep	4.938	536	.003	.488			

Substituting from above, we obtain, for our estimated odds, the following.

Model	M and S (condition on W)	M and W (condition on S)	S and W (condition on M)
Saturated – level 0	$\widehat{\Theta_{MS(0)}} = e^{1.36 + 0.088} = 4.26$	$\widehat{\Theta_{M(0)W}} = e^{.694 + .088} = 2.19$	$\widehat{\Theta_{(0)SW}} = e^{.352 + .088} =$
			1.55
Saturated – level 1	$\widehat{\Theta_{MS(1)}} = e^{1.36} = 3.90$	$\widehat{\Theta_{M(1)W}} = e^{.694} = 2.00$	$\widehat{\Theta_{(1)SW}} = e^{.352} = 1.42$
Homogenous	$\widehat{\Theta_{MS(0)}} = \widehat{\Theta_{MS(1)}} = e^{1.396} =$	$\widehat{\Theta_{M(0)W}} = \widehat{\Theta_{M(1)W}} = e^{.748} =$	$\widehat{\Theta_{(0)SW}} = \widehat{\Theta_{(1)SW}} = e^{.385} =$
	4.04	2.11	1.47
Cond. Indep (SW)	$\widehat{\Theta_{MS(0)}} = \widehat{\Theta_{MS(1)}} = e^{1.464} =$	$\widehat{\Theta_{M(0)W^{\pm}}}  \widehat{\Theta_{M(1)W}} = e^{.875} =$	1
	4.32	2.40	
Independence	1	1	1

We note that we can also check these numbers by calculating the estimated log(odds) using the appropriate 2x2 table of expected frequencies. We include the following summary table of observed and expected frequencies for our models (these numbers were obtained from the cell counts and residuals tables in our output).

					Expe	ected	
Manage	Superv	Worker	Observed	Saturated	Homogenous	Condl Indep	Independ
						Superv-Wrkr	
Bad(0)	Low(0)	Low(0)	103.00	103.00	102.26	97.16	50.29
Bad(0).	Low(0)	High(1)	87.00	87.00	87.74	92.84	81.90
Bad(0)	High(1)	Low(0)	32.00	32.00	32.74	37.84	50.15
Bad(0)	High(1)	High(1)	42.00	42.00	41.26	36.16	81.67
High(1)	Low(0)	Low(0)	59.00	59.00	59.74	51.03	85.91
High(1)	Low(0)	High(1)	109.00	109.00	108.26	116.90	139.91
High(1)	High(1)	Low(0)	78.00	78.00	77.26	85.97	85.67
High(1)	High(1)	High(1)	205.00	205.00	205.74	197.03	139.52

Mod	M and S (condition on W)	M and W (condition on S)	S and W (condition on M)
Sat-0	(103x78)/(32x59)=4.26	(103x109)/(87x59) = 2.19	(103x42)/(32x87) = 1.554
Sat – 1	(87x205)/(42x109)=3.90	(32x205)/(42x78)=2.00	(59x205)/(78x109) = 1.423
Homo	(102.26x77.26)/(32.74x59.74)=4.04	(102.26x108.26)/(87.74x59.74)=2.11	(102.26x41.26)/(87.74x32.74)=1.47
	(87.74x205.74)/(41.26x108.26)=4.04	(32.74x205.74)/(41.26x77.26)=2.11	(59.74x205.75)/(108.26x77.26)=1.47
Cond In(SW)	(97.16x85.97)/(37.84x51.03)=4.32	(97.16x116.9)/(92.84x51.03) = 2.40	(97.16x36.16)/(92.84x37.84)=1.00
	(92.84x197.03)/(36.16x116.90)=4.32	(37.84x197.03)/(36.16x85.97)=2.40	(51.03x197.03)/116.90x85.97)= 1.00
Indep	(50.29x85.67)/(50.15x85.91) = 1.00	(50.29x139.91)/(81.90x85.91)=1.00	(50.29x81.67)/(81.90x50.15)=1.00
	(81.90x139.52)/(81.67x139.91)=1.00	(50.15x139.52)/(81.67x85.67)=1.00	(85.91x139.52)/(139.91x85.67)=1.00

We illustrate how the entry is obtained for one of the cells above. Other entries are obtained in an analogous manner.

Consider  $\widehat{\Theta_{M(0)W}}$  for our Homogenous Model.

For S =0, our  $2x^2$  table of expected frequencies follows.

	W(0)	W(1)
M(0)	102.26	87.74
M(1)	59.74	108.26

 $\widehat{\Theta_{MS(0)}} = (103)(78)/(59)(32)$