

# Logistic Regression with Polytomous & Ordinal Data

EDMS 771

Brandi A. Weiss

1

## A Look Back & A Look Ahead

### □ Last Week:

#### ■ Binary Logistic Regression

- DV = dichotomous

### □ This Week:

#### ■ Polytomous Logistic Regression

- aka Multinomial LR
- DV = 3 or more nominal categories

#### ■ Ordinal Logistic Regression

- DV = 2 or more ordered categories



2

## Review: Coin Flip Example



- Probability coin will be heads:
  - 50%
- Odds coin will be heads:
  - $\frac{\text{probability coin will be heads}}{\text{probability coin will be tails}} = \frac{50}{50} = 1$
- Logit coin will be heads:
  - $\ln(\text{odds}) = \ln(50/50) = 0$

3

## LR with Polytomous Data (aka Multinomial LR)

5

## Running Example #1:

---

- ☐ **Research Question:** Do beliefs about impact of working mothers on children aid in predicting a women's work status?
- ☐ **Data: GSS** (<http://sda.berkeley.edu/index.htm>)
  - Only data from females used
- ☐ **Independent Variable (1-4 scale):**
  - A working mother can establish just as warm and secure a relationship with her children as a mother who does not work. (*fechld*)
    - ☐ Strongly Agree
    - ☐ Agree
    - ☐ Disagree
    - ☐ Strongly Disagree

6

## Running Example #1 (continued)

---

### **Binary LR:**

- ☐ **DV:**
  - Working (interest group)
  - Not Working (reference group) (i.e., Unemployed)

### **Polytomous LR:**

- ☐ **DV:**
  - Working (interest group)
  - Retired (interest group)
  - Student (interest group)
  - Not Working (reference group)

7

## Why Not Run 4 Binary LRs???



- Could run 3 Binary LRs:
  - Working vs Not Working
  - Retired vs Not Working
  - Student vs Not Working
- Problem: Lose a lot of info!
  - Example: if we use Student vs Non-Student, we may NOT get a statistically significant effect for one IV (e.g., gender) if males:
    - Are less likely to be Unemployed
    - And more likely to be Working

8

## Polytomous LR

- Odds = interest group/reference group
- Need k-1 regression equations (i.e., one for each interest group):
  - Logit (working) =  $a_w + b_w X$
  - Logit (retired) =  $a_r + b_r X$
  - Logit (student) =  $a_s + b_s X$
- Classify/predict people into one of 4 categories

*Note:* For ease of reading the notes, for this lecture

$$\exp[a + bX] = e^{a + bX}$$

9

# SPSS

**[DataSet2] - SPSS Data Editor**

File Edit View Analyze Reports Graphs Utilities Add-ons Window Help

1  
rkstat  
7  
7  
5  
7  
1  
5  
2  
1  
7  
7  
1

age marital fechld  
54 1  
71 2

Regression  
Loglinear  
Classify  
Data Reduction  
Scale  
Nonparametric Tests  
Survival  
Multiple Response

Linear...  
Curve Estimation...  
Binary Logistic...  
**Multinomial Logistic...**  
Ordinal...  
Probit...  
Nonlinear...  
Weight Estimation...  
2-Stage Least Squares...

**Multinomial Logistic Regression**

Dependent: work4(Last)

Reference Category...

Factor(s):

Covariate(s): MOTHER WORKING

Model... Statistics... Criteria... Options... Save...

10

# SPSS

**Multinomial Logistic Regression...**

Reference Category

☐ First category  
☒ Last category  
☐ Custom

Value:

Category Order

☒ Ascending  
☐ Descending

Continue  
Cancel  
Help

**Multinomial Logistic Regression: Statistics**

☒ Case processing summary

Model

☒ Pseudo R-square  
☒ Step summary  
☒ Model fitting information  
☐ Information criteria

☐ Cell probabilities  
☒ Classification table  
☐ Goodness-of-fit

Parameters

☒ Estimates  
☒ Likelihood ratio tests  
☐ Asymptotic correlations  
☐ Asymptotic covariances

Confidence Interval (%): 95

Define Subpopulations

☒ Covariate patterns defined by factors and covariates  
☐ Covariate patterns defined by variable list below

Subpopulations:

fechld( )

Continue  
Cancel  
Help

In this example DV was coded as follows:  
1=Working  
2=Retired  
3=Student  
4=Not Working  
Recall we said "Not Working" would be the reference group

## SPSS Output

Case Processing Summary			
		N	Marginal Percentage
work4	Working	5679	55.4%
	Retired	1221	11.9%
	Student	292	2.9%
	Not Working	3052	29.8%
Valid		10244	100.0%
Missing		14298	
Total		24542	
Subpopulation		4	

With a Blind Model, if you predicted everyone to be "Working", you'd be correct 55.4% of the time.

12

## SPSS Output (continued)

Pseudo R-Square	
Cox and Snell	.059
Nagelkerke	.067
McFadden	.029

- ☐ Like  $R^2$  values
- ☐ Not recommended to use
- ☐ Cox & Snell will never = 1.0
- ☐ Nagelkerke scaled Cox & Snell's computation so that it will go to 1.0
- ☐ McFadden = Likelihood Ratio  $R^2$

13

## SPSS Output

The 2 models are nested, therefore:

$$\chi^2_{\text{difference}} = 744.048 - 123.544 = 620.536$$

with 3 df,  $p < .001$

Therefore, the Full model (with **fechld** as a predictor) is significantly better than the 'blind' model.

Model Fitting Information				
Model	Model Fitting Criteria	Likelihood Ratio Tests		
	-2 Log Likelihood	Chi-Square	df	Sig.
Intercept Only	744.048			
Final	123.544	620.504	3	.000

1 df for each regression equation. We have 4 groups, and  $k-1$  regression equations.  
 $4 - 1 = 3$  df

14

## Full/Final Model vs Reduced Model

- **Final Model** – includes the predictor:
  - 'Fechld'
- **Reduced Model** – removes the predictor
  - Ex: 'Fechld' Reduced Model contains all predictors except 'Fechld'
    - In this case: Reduced model = Blind model, because only 1 predictor
    - ∴ redundant with "Model Fitting Information" table in SPSS
  - Likelihood Ratio Test:
    - $\chi^2 = -2LL_{\text{REDUCED}} - 2LL_{\text{FINAL}}$
    - allows you to test contribution of 'Fechld'

15

## SPSS Output (continued)

The 2 models are nested, therefore:

$$\chi^2_{\text{difference}} = 744.048 - 123.544 = \mathbf{620.536} \text{ with 3 df, } p < .001$$

Therefore, the Full model (with **fechld** as a predictor) is significantly better than the 'blind' model.

### Likelihood Ratio Tests

Model Fitting Criteria				
Likelihood Ratio Tests				
Effect	-2 Log Likelihood of Reduced Model	Chi-Square	df	Sig.
Intercept	2037.841	1914.297	3	.000
fechld	744.048	620.504	3	.000

The chi-square statistic is the difference in -2 log-likelihoods between the final model and a reduced model. The reduced model is formed by omitting an effect from the final model. The null hypothesis is that all parameters of that effect are 0.

### Likelihood Ratio Test:

$$\chi^2 = -2LL_{\text{REDUCED}} - 2LL_{\text{FINAL}}$$

Tests if variable "fechld" contributes statistically significantly to the model.

(a little redundant w/ the 1 predictor model)

16

## SPSS Output (continued)

### Parameter Estimates

							95% Confidence Interval for Exp(B)		
work4 <sup>a</sup>		B	Std. Error	Wald	df	Sig.	Exp(B)	Lower Bound	Upper Bound
Working	Intercept	1.821	.062	872.903	1	.000			
	fechld	-.566	.026	461.876	1	.000	.568	.539	.598
Retired	Intercept	-.976	.094	108.062	1	.000			
	fechld	.025	.037	.465	1	.495	1.026	.954	1.103
Student	Intercept	-1.065	.156	46.726	1	.000			
	fechld	-.609	.075	66.684	1	.000	.544	.470	.630

a. The reference category is: Not Working.

- $\text{Logit}_w = 1.821 - .566(\text{fechld})$
- $\text{Logit}_r = -0.976 + .025(\text{fechld})$
- $\text{Logit}_s = -1.065 - .609(\text{fechld})$

17



## SPSS Output: Parameter Estimates (continued)

- **Working:** If IV changed 1 unit, the DV-odds would change by a multiplicative amount of .568
- Odds of being Student to Not Working are much lower than Working to Not Working
- Look at predictors for each regression equation separately!
  - **Belief about working mothers is not a significant predictor of whether or not a person is retired or Not Working**

18

## SPSS Output (continued)



Classification					
Observed	Predicted				Percent Correct
	Working	Retired	Student	Not Working	
Working	5415	0	0	264	95.4%
Retired	1096	0	0	125	.0%
Student	281	0	0	11	.0%
Not Working	2700	0	0	352	11.5%
Overall Percentage	92.7%	.0%	.0%	7.3%	56.3%

- Only classifying 56.3% of women correctly
- Not classifying anyone as Retired or Student
  - Hmm, small percentage of respondents in these categories could → problems
- Blind model predicted women work status correctly 55.4%

22

## Classification Rule: Binary LR

---

□ 2 equations:

- $P(\text{Working}) = \frac{\exp[a + bX]}{1 + \exp[a + bX]}$

- $P(\text{Not-Working}) = \frac{1}{1 + \exp[a + bX]}$

23

## Classification Rule: Polytomous LR

---

Recall:

Odds =  $\frac{\text{odds of being in a particular group}}{\text{sum of odds of being in all groups}}$

□ (Relative to reference group)

□ 4 equations:

- $P(\text{Group } j) = \frac{\exp[a_j + b_j X]}{1 + \sum \exp[a_j + b_j X]}$

24

## Classification Rule: Polytomous LR

□ 4 equations:

$$\blacksquare P(\text{Working}) = \frac{\exp[a_W + b_W X]}{1 + \exp[a_W + b_W X] + \exp[a_R + b_R X] + \exp[a_S + b_S X]}$$

$$\blacksquare P(\text{Retired}) = \frac{\exp[a_R + b_R X]}{1 + \exp[a_W + b_W X] + \exp[a_R + b_R X] + \exp[a_S + b_S X]}$$

$$\blacksquare P(\text{Student}) = \frac{\exp[a_S + b_S X]}{1 + \exp[a_W + b_W X] + \exp[a_R + b_R X] + \exp[a_S + b_S X]}$$

$$\blacksquare P(\text{Not-Working}) = \frac{1}{1 + \exp[a_W + b_W X] + \exp[a_R + b_R X] + \exp[a_S + b_S X]}$$

25

## Regression Equations:

□ If  $\text{Fechld} = 4$  (i.e., believe mother working negatively impacts children) then

$$\blacksquare \text{Logit}_{\text{working}} = 1.821 - .566(4) = -.443$$

$$\blacksquare \text{Logit}_{\text{retired}} = -0.976 + .025(4) = -.876$$

$$\blacksquare \text{Logit}_{\text{student}} = -1.065 - .609(4) = -3.501$$

$$\blacksquare \exp(\text{Logit}_{\text{working}}) = \exp(-.443) = .642$$

odds( Working to Not Working)

$$\blacksquare \exp(\text{Logit}_{\text{retired}}) = \exp(-.876) = .416$$

odds( Retired to Not Working)

$$\blacksquare \exp(\text{Logit}_{\text{student}}) = \exp(-3.501) = .030$$

odds( Student to Not Working)

26

## Regression Equations:

---

- If  $Fechld = 4$  (i.e., believe mother working negatively impacts children) then

- $P(\text{Working}) = \frac{.642}{1 + .642 + .416 + .030} = .31$

- $P(\text{Retired}) = \frac{.416}{1 + .642 + .416 + .030} = .20$

- $P(\text{Student}) = \frac{.030}{1 + .642 + .416 + .030} = .01$

- $P(\text{Not Working}) = \frac{1}{1 + .642 + .416 + .030} = .48$

27

## Running Example #2

---

### Polytomous LR w/ Multiple Predictors:

- DV:
  - Working (interest group)
  - Retired (interest group)
  - Student (interest group)
  - Not Working (reference group)
- IV's:
  - $Fechld$  – Belief about working mothers
  - Years of Education
  - Age
  - Marital Status (dummy coded – 1=married)

28

# SPSS

**[DataSet2] - SPSS Data Editor**

File Edit View Analyze Reports Graphs Utilities Add-ons Window Help

1  
rkstat  
7  
7  
5  
7  
7  
5  
2  
1  
7  
7  
7  
1  
2  
13

age marital fechld  
54 1  
71 2

Regression  
Loglinear  
Classify  
Data Reduction  
Scale  
Nonparametric Tests  
Survival  
Multiple Response

Linear...  
Curve Estimation...  
Binary Logistic...  
**Multinomial Logistic...**  
Ordinal...  
Probit...  
Nonlinear...  
Weight Estimation...  
2-Stage Least Squares...

**Multinomial Logistic Regression**

Dependent:  
work4(Last)

Reference Category...

Factor(s):

Covariate(s):  
MOTHER WORK  
HIGHEST YEAR C  
AGE OF RESPON  
marrydummy

Model... Statistics... Criteria... Options... Save...

29

# SPSS

**Multinomial Logistic Regression: Statistics**

☒ Case processing summary

Model  
☒ Pseudo R-square  
☒ Step summary  
☒ Model fitting information  
☐ Information criteria

☐ Cell probabilities  
☒ Classification table  
☐ Goodness-of-fit

Parameters  
☒ Estimates  
☒ Likelihood ratio tests  
☐ Asymptotic correlations  
☐ Asymptotic covariances

Confidence Interval (%): 95

Define Subpopulations  
☒ Covariate patterns defined by factors and covariates  
☐ Covariate patterns defined by variable list below

Subpopulations:  
fechld( )

In this example DV was coded as follows:  
1=Working  
2=Retired  
3=Student  
4=Not Working  
Recall we said "Not Working" would be the reference group

## SPSS Output

Case Processing Summary

		N	Marginal Percentage
work4	Working	5645	55.5%
	Retired	1209	11.9%
	Student	290	2.8%
	Not Working	3035	29.8%
Valid		10179	100.0%
Missing		14363	
Total		24542	
Subpopulation		3734 <sup>a</sup>	

a. The dependent variable has only one value observed in 2500 (67.0%) subpopulations.

Pseudo R-Square

Cox and Snell	.423
Nagelkerke	.483
McFadden	.264

31

## SPSS Output

The 2 models are nested, therefore:

$$\chi^2_{\text{difference}} = 15424.231 - 9823.910 = \mathbf{5600.321} \text{ with 12 df, } p < .001$$

Therefore, the Full model (with **fechld**, **education**, **age**, and **marital status** as a predictors) is significantly better than the 'blind' model.

Model Fitting Information

Model	Model Fitting Criteria	Likelihood Ratio Tests		
	-2 Log Likelihood	Chi-Square	df	Sig.
Intercept Only	15424.231			
Final	9823.910	5600.321	12	.000

1 df for each predictor in each regression equation. We have 4 predictors, 4 groups, and k-1 regression equations.  
 $4 - 1 = 3$  regression equations \* 4 predictors = 12

32

## Full/Final Model vs Reduced Model

- **Final Model** - has all 4 predictors:
  - 'Fechld'
  - Education
  - Age
  - Marital Status
- **Reduced Model** – removes 1 of the predictors
  - Ex: 'Fechld' Reduced Model contains all predictors except 'Fechld'
  - Likelihood Ratio Test:
    - $\chi^2 = -2LL_{\text{REDUCED}} - 2LL_{\text{FINAL}}$
    - allows you to test contribution of 'Fechld'

33

## SPSS Output (continued)

Does the variable contribute statistically significantly to the model?

$$\text{Likelihood Ratio Test} = 10012.569 - 9823.910 = 188.659$$

Likelihood Ratio Tests				
Model Fitting Criteria		Likelihood Ratio Tests		
Effect	-2 Log Likelihood of Reduced Model	Chi-Square	df	Sig.
Intercept	11079.234	1255.323	3	.000
fechld	10012.569	188.659	3	.000
educ	10420.064	596.154	3	.000
age	13443.256	3619.346	3	.000
marrydummy	10031.433	207.523	3	.000

The chi-square statistic is the difference in -2 log-likelihoods between the final model and a reduced model. The reduced model is formed by omitting an effect from the final model. The null hypothesis is that all parameters of that effect are 0.

The Final model fits significantly better than a model WITHOUT 'fechld' in it (i.e., 'fechld' significantly contributes to the model).

All 4 predictors significantly contribute to the model

34

## SPSS Output (continued)

Parameter Estimates								
		B	Std. Error	Wald	df	Sig.	95% Confidence Interval for Exp(B)	
work4 <sup>a</sup>							Lower Bound	Upper Bound
Working	Intercept	-.850	.186	20.928	1	.000		
	fechld	-.375	.028	176.360	1	.000	.687	.650 .726
	educ	.221	.010	492.222	1	.000	1.247	1.223 1.271
	age	-.031	.002	336.890	1	.000	.970	.967 .973
	marrydummy	.563	.050	124.692	1	.000	1.755	1.590 1.938
Retired	Intercept	-10.897	.379	828.497	1	.000		
	fechld	-.309	.049	39.502	1	.000	.734	.667 .809
	educ	.161	.014	126.106	1	.000	1.175	1.142 1.209
	age	.125	.004	957.222	1	.000	1.133	1.124 1.142
	marrydummy	.661	.090	54.274	1	.000	1.937	1.625 2.310
Student	Intercept	-2.625	.556	22.333	1	.000		
	fechld	-.339	.081	17.451	1	.000	.712	.608 .835
	educ	.275	.028	97.125	1	.000	1.317	1.246 1.391
	age	-.141	.009	263.167	1	.000	.868	.853 .883
	marrydummy	1.521	.157	93.666	1	.000	4.575	3.362 6.224

a. The reference category is: Not Working.

## Parameter Estimates → Regression Equations

- $\text{Logit}_w = -.850 - .375(\text{fechld}) + .221(\text{educ}) - .031(\text{age}) + .563(\text{marital})$
- $\text{Logit}_R = -10.897 - .309(\text{fechld}) + .161(\text{educ}) + .125(\text{age}) + .661(\text{marital})$
- $\text{Logit}_S = -2.625 - .339(\text{fechld}) + .275(\text{educ}) - .141(\text{age}) + 1.521(\text{marital})$
- **Working:**
- If 'fechld' changed 1 unit, the log-odds of the DV would change by a multiplicative amount of .687, all else being held constant
- If 'educ' changed 1 unit, the log-odds of the DV would change by a multiplicative amount of 1.247, all else being held constant



## SPSS Output (continued)

Classification					
Observed	Predicted				Percent Correct
	Working	Retired	Student	Not Working	
Working	4975	68	0	602	88.1%
Retired	175	718	0	316	59.4%
Student	284	0	0	6	.0%
Not Working	1690	411	0	934	30.8%
Overall Percentage	70.0%	11.8%	.0%	18.3%	65.1%

- Now classifying 65.1% of women correctly
- Not classifying anyone as Student
- Blind model predicted women work status correctly 55.4%
- 1 predictor model predicted women correctly 56.3%

40

## LR with Ordinal Data

41

## Measurement Scales

---

- **Nominal:**
  - Categories; no order
  - Ex: Goal/No Goal
- **Ordinal:**
  - Rank order but cannot measure "by how much"
  - Ex: Degree (HS, Associate, Bachelor, Master, Doctorate)
- **Interval:**
  - Rank order & equidistance
  - Ex: Years of Education

42

## My DV is Ordinal; What Should I Do??

---

- Pretend variable is on interval scale → OLS
- Treat as though it were *measured* on ordinal scale but is *really* interval/ratio underneath → WLS
- **\*\*Treat as though it was measured on a true ordinal scale → Cumulative logit model**
- Treat it as nominal → Polytomous LR

Slide Adapted from Roy Levy

43

### Running Example #3:

---

- **Research Question:** Does father's education (faeduc) aid in predicting a person's highest degree obtained?
- **Data:** GSS (<http://sda.berkeley.edu/index.htm>)
- **DV = Highest Degree Obtained:**
  - Dropout (interest group)
  - High School (interest group)
  - Some College (interest group)
  - Bachelor (interest group)
  - Graduate (interest group)
- **Independent Variable =** years of father's education

44

### Now Predicting Cumulative Probabilities:

---

Drop	vs.	All other classes (HS, Some, Bachelor, Graduate)
Drop & HS	vs.	Some, Bachelor, Graduate
Drop, HS, & Some	vs.	Bachelor & Graduate
Drop, HS, Some, Bachelor	vs.	Graduate

- Need k-1 regression equations ( $5 - 1 = 4$ )

45

## "Cumulative" Logits

- **Binary:**
  - Logit expressions = probabilities of being in a given group
- **Polytomous:**
  - Logit expressions = probability of being in a given group compared to probability of being in reference group
- **Ordinal:**
  - Logit expressions = probabilities of being *in or below* a given group
    - i.e., "cumulative" logits
  - Logit ( $\leq$ Drop) = logit for being *in or below* Dropout
  - Logit ( $\leq$ HS) = logit for being *in or below* HS
  - Logit ( $\leq$ SC) = logit for being *in or below* SC
  - Logit ( $\leq$ Bachelor) = logit for being *in or below* Bachelor<sup>46</sup>

## Homogeneity of Regression

We know that:

- $\text{Logit}(\leq \text{Drop}) < \text{Logit}(\leq \text{HS}) < \text{Logit}(\leq \text{SC}) < \text{Logit}(\leq \text{Bachelor})$

→ Assumption:

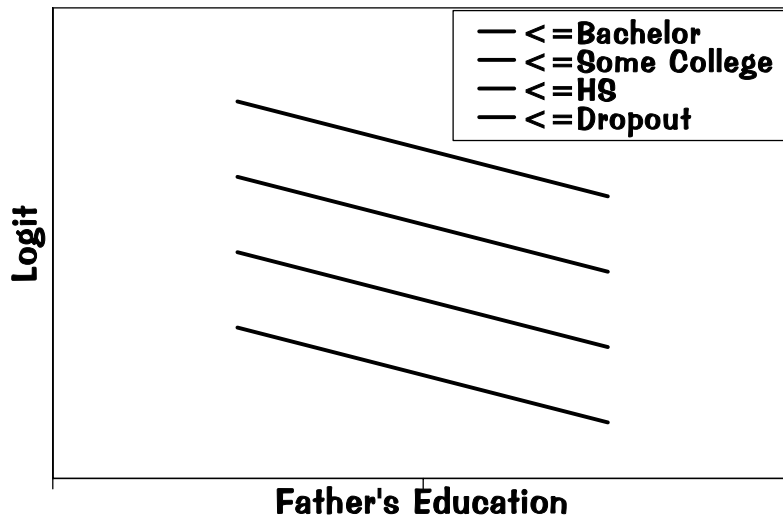
- $\text{Logit}(\leq \text{Drop}) = a_D - b_1 X_1 - b_2 X_2 - \dots - b_p X_p$
- $\text{Logit}(\leq \text{HS}) = a_H - b_1 X_1 - b_2 X_2 - \dots - b_p X_p$
- $\text{Logit}(\leq \text{SC}) = a_C - b_1 X_1 - b_2 X_2 - \dots - b_p X_p$
- $\text{Logit}(\leq \text{Bachelor}) = a_B - b_1 X_1 - b_2 X_2 - \dots - b_p X_p$

- Slopes ( $b_1, b_2, \dots, b_p$ ) are the **SAME** for each equation
- Only the intercepts ( $a_D, a_H, a_C, a_B$ ) differ

47

# Homogeneity of Regression

## □ Parallel Lines



48

## SPSS

SPSS Data Editor - .sav [DataSet1] - SPSS Data Editor

Menu: Transform, Analyze, Graphs, Utilities, Add-ons, Window, Help

Submenu: Reports, Descriptive Statistics, Compare Means, General Linear Model, Mixed Models, Correlate, Regression, Loglinear, Classify, Data Reduction, Scale, Nonparametric Tests, Survival, Multiple Response

Regression Submenu: Linear..., Curve Estimation..., Binary Logistic..., Multinomial Logistic..., Ordinal..., Probit..., Nonlinear..., Weight Estimation..., 2-Stage Least Squares...

Ordinal Regression Dialog Box:

- Dependent: RS HIGHEST DEGREE
- Factor(s):
- Covariate(s): TOTAL FAMILY INCOME
- Case Identification Variables: RESPONDENTS SEX, AGE OF RESPONDER, MOTHERS HIGHEST, FATHERS HIGHEST, SPOUSES HIGHEST, MOTHER WORKING
- Buttons: OK, Paste, Reset, Cancel, Help, Options..., Output..., Location..., Scale...

49

## SPSS (continued)

**Ordinal Regression: Output**

Display

☐ Print iteration history for every 1 step(s)

☒ Goodness of fit statistics

☒ Summary statistics

☒ Parameter estimates

☐ Asymptotic correlation of parameter estimates

☐ Asymptotic covariance of parameter estimates

☐ Cell information

☒ Test of parallel lines

Saved variables

☒ Estimated response probabilities

☒ Predicted category

☐ Predicted category probability

☐ Actual category probability

Print log-likelihood

☒ Including multinomial constant

☐ Excluding multinomial constant

Continue Cancel Help

50

## SPSS Output: Sparse Data!

### Warnings

There are 1 (1.0%) cells (i.e., dependent variable levels by combinations of predictor variable values) with zero frequencies.

- ☐ We haven't seen this message before!! ☹
- ☐ Problem: There weren't any respondents who had father's with 1 year of education who were also in the "some college" category.
- ☐ 1 Solution: Gather more data
- ☐ There was only 1 cell in this dataset with sparse data: we'll ignore for now

		DV (Highest Degree Obtained)				
		Dropout	HS	Some College	Bachelor	Grad
Father's Education (years)	1	52	29	0	2	2
	2	152	97	3	14	8
	...					
	20	13	202	33	232	148

51

## SPSS Output

Case Processing Summary			
		N	Marginal Percentage
RS HIGHEST DEGREE	DROPOUT	5733	18.3%
	HIGH SCHOOL	16872	53.9%
	SOME COLLEGE	1610	5.1%
	BACHELOR	4842	15.5%
	GRADUATE	2241	7.2%
Valid		31298	100.0%
Missing		12400	
Total		43698	

Pseudo R-Square	
Cox and Snell	.195
Nagelkerke	.211
McFadden	.085
Link function: Logit.	

52

## SPSS Output (continued)

Test of Parallel Lines <sup>a</sup>				
Model	-2 Log Likelihood	Chi-Square	df	Sig.
Null Hypothesis	1025.101			
General	810.696	214.404	3	.000

The null hypothesis states that the location parameters (slope coefficients) are the same across response categories.

a. Link function: Logit.

- **Tests Homogeneity of Regression Assumption**
  - **Reject  $H_0$**  – The slopes are different across logits ☹️
  - **Often anti-conservative** – i.e., reports p-values that are too low with large samples or continuous predictors
  - **We'll continue anyway for illustration purposes**

53

## SPSS Output (continued)

$$\chi^2 = 7801.279 - 1025.101 = 6776.179$$

with 1 df,  $p < .001$

Therefore, the Full model (with **faeduc** as a predictor) is significantly better than the 'blind' model.

### Model Fitting Information

Model	-2 Log Likelihood	Chi-Square	df	Sig.
Intercept Only	7801.279			
Final	1025.101	6776.179	1	.000

Link function: Logit.

### Goodness-of-Fit

	Chi-Square	df	Sig.
Pearson	591.180	79	.000
Deviance	534.314	79	.000

Link function: Logit.

All texts define these slightly differently. In general tests if your model fits well compared to a saturated model. Let's ignore this for now.

54

## SPSS Output (continued)

### Parameter Estimates

		Estimate	Std. Error	Wald	df	Sig.	95% Confidence Interval	
							Lower Bound	Upper Bound
Threshold	[degree = 0]	.546	.029	363.888	1	.000	.490	.602
	[degree = 1]	3.448	.035	9443.233	1	.000	3.379	3.518
	[degree = 2]	3.762	.036	10682.037	1	.000	3.691	3.833
	[degree = 3]	5.222	.042	15201.079	1	.000	5.139	5.305
Location	paeduc	.222	.003	6198.658	1	.000	.217	.228

Link function: Logit.

- **SPSS: Logit = a - bX**
  - Note the Minus sign!!

HRA: Notice the Same Slope for all equations

- **Logit (≤ Dropout) = 0.546 - .222 (faedu)**
- **Logit (≤ HS) = 3.448 - .222 (faedu)**
- **Logit (≤ SC) = 3.762 - .222 (faedu)**
- **Logit (≤ Bachelor) = 5.222 - .222 (faedu)**

55



## Probabilities & Category Predictors

\*GSS Degree Dataset.sav [DataSet2] - SPSS Data Editor

File Edit View Data Transform Analyze Graphs Utilities Add-ons Window Help

1 : CASEID 20020001

	EST1_1	EST2_1	EST3_1	EST4_1	EST5_1	PRE_1
22953	.11	.58	.06	.18	.07	1
22954	.23	.62	.04	.09	.03	1
22955	.23	.62	.04	.09	.03	1
22956	.13	.60	.06	.15	.06	1
22957	.47	.47	.02	.03	.01	1
22958	.	.	.	.	.	.
22959	.31	.58	.03	.06	.02	1
22960	.31	.58	.03	.06	.02	1
22961	.31	.58	.03	.06	.02	1
22962	.23	.62	.04	.09	.03	1
22963	.11	.58	.06	.18	.07	1
22964	.	.	.	.	.	.
22965	.53	.43	.01	.03	.01	0
22966	.53	.43	.01	.03	.01	0
22967	.11	.58	.06	.18	.07	1
22968	.19	.62	.04	.11	.04	1
22969	.11	.58	.06	.18	.07	1
22970	.11	.58	.06	.18	.07	1
22971	.13	.60	.06	.15	.06	1
22972	.19	.62	.04	.11	.04	1
22973	.23	.62	.04	.09	.03	1
22974	.16	.62	.05	.13	.05	1
22975	.	.	.	.	.	.
22976	.13	.60	.06	.15	.06	1
22977	.23	.62	.04	.09	.03	1
22978	.	.	.	.	.	.

59

## What are our Cut-Points?

<u>years faeduc</u>	<u>p(dropout)</u>	<u>p(hs)</u>	<u>p(sc)</u>	<u>p(bachelor)</u>	<u>p(grad)</u>
0	<b>0.63</b>	0.34	0.01	0.02	0.01
1	<b>0.58</b>	0.38	0.01	0.02	0.01
2	<b>0.53</b>	0.43	0.01	0.03	0.01
3	<b>0.469</b>	0.472	0.02	0.03	0.01
4	0.42	<b>0.51</b>	0.02	0.04	0.01
5	0.36	<b>0.55</b>	0.02	0.05	0.02
6	0.31	<b>0.58</b>	0.03	0.06	0.02
7	0.27	<b>0.60</b>	0.03	0.07	0.02
8	0.23	<b>0.62</b>	0.04	0.09	0.03
9	0.19	<b>0.62</b>	0.04	0.11	0.04
10	0.16	<b>0.62</b>	0.05	0.13	0.05

60

## What are our Cut-Points?

<u>years faeduc</u>	<u>p(dropout)</u>	<u>p(hs)</u>	<u>p(sc)</u>	<u>p(bachelor)</u>	<u>p(grad)</u>
11	0.13	<b>0.60</b>	0.06	0.15	0.06
12	0.11	<b>0.58</b>	0.06	0.18	0.07
13	0.09	<b>0.55</b>	0.07	0.21	0.09
14	0.07	<b>0.51</b>	0.07	0.23	0.11
15	0.06	<b>0.47</b>	0.08	0.26	0.13
16	0.05	<b>0.43</b>	0.08	0.29	0.16
17	0.04	<b>0.38</b>	0.08	0.31	0.19
18	0.03	<b>0.335</b>	0.07	0.332	0.22
19	0.02	0.29	0.07	<b>0.34</b>	0.27
20	0.02	0.25	0.07	<b>0.35</b>	0.32

61

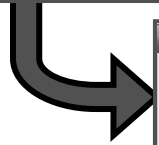
## Running Example #4

### Ordinal LR w/ Multiple Predictors:

- DV = Highest Degree Obtained:
  - Dropout (interest group)
  - High School (interest group)
  - Some College (interest group)
  - Bachelor (interest group)
  - Graduate (interest group)
- IV's:
  - Father's Education in years (faeduc)
  - Mother's Education in years (maeduc)
  - Gender

62

The screenshot shows the SPSS Data Editor window with the 'Analyze' menu open. The 'Regression' submenu is also open, and the 'Ordinal...' option is highlighted. The background data table is partially visible, showing variables 'free', 'nadeg', 'padeg', and 'spdeg'.



**Ordinal Regression**

Case Identification

- AGE OF RESPON
- MOTHERS HIGH
- FATHERS HIGH
- SPOUSES HIGH
- HIGHEST YEAR
- MOTHER WORKI
- RESPONDENTS
- TOTAL FAMILY I
- Estimated Cell Pro
- Estimated Cell Pro
- Estimated Cell Pro
- Estimated Cell Pro

Dependent:

RS HIGHEST DEGREE

Factor(s)

Covariate(s)

- HIGHEST YEAR SCH
- HIGHEST YEAR SCH
- RESPONDENTS SEX

Options... Output... Location... Scale...

## SPSS (continued)

**Ordinal Regression: Output**

Display

- ☐ Print iteration history for every  step(s)
- ☒ Goodness of fit statistics
- ☒ Summary statistics
- ☒ Parameter estimates
- ☐ Asymptotic correlation of parameter estimates
- ☐ Asymptotic covariance of parameter estimates
- ☐ Cell information
- ☒ Test of parallel lines

Saved variables

- ☒ Estimated response probabilities
- ☒ Predicted category
- ☐ Predicted category probability
- ☐ Actual category probability

Print log-likelihood

- ☒ Including multinomial constant
- ☐ Excluding multinomial constant

Continue Cancel Help

27

## SPSS Output

### Warnings

There are 1316 (38.4%) cells (i.e., dependent variable levels by combinations of predictor variable values) with zero frequencies.

### Case Processing Summary

		N	Marginal Percentage
RS HIGHEST DEGREE	DROPOUT	4932	16.9%
	HIGH SCHOOL	15787	54.2%
	SOME COLLEGE	1540	5.3%
	BACHELOR	4695	16.1%
	GRADUATE	2181	7.5%
Valid		29135	100.0%
Missing		14563	
Total		43698	

65

## SPSS Output (continued)

### Model Fitting Information

Model	-2 Log Likelihood	Chi-Square	df	Sig.
Intercept Only	15239.670			
Final	7685.279	7554.391	3	.000

Link function: Logit.

### Goodness-of-Fit

	Chi-Square	df	Sig.
Pearson	4569.821	2737	.000
Deviance	3543.083	2737	.000

Link function: Logit.

### Pseudo R-Square

Cox and Snell	.228
Nagelkerke	.248
McFadden	.102

Link function: Logit.

66

## SPSS Output (continued)

Parameter Estimates								
		Estimate	Std. Error	Wald	df	Sig.	95% Confidence Interval	
							Lower Bound	Upper Bound
Threshold	[degree = 0]	1.015	.052	379.292	1	.000	.913	1.117
	[degree = 1]	4.068	.058	4943.722	1	.000	3.955	4.181
	[degree = 2]	4.388	.059	5624.599	1	.000	4.273	4.503
	[degree = 3]	5.873	.063	8747.770	1	.000	5.750	5.996
Location	paeduc	.139	.004	1403.286	1	.000	.132	.146
	maeduc	.158	.004	1277.636	1	.000	.149	.167
	sex	-.127	.023	30.564	1	.000	-.172	-.082

Link function: Logit.

- ☐ Hmmm, No likelihood-ratio tests ☹️
- ☐ Need to compute them yourself

67

## SPSS Output (continued)

Test of Parallel Lines <sup>a</sup>				
Model	-2 Log Likelihood	Chi-Square	df	Sig.
Null Hypothesis	7685.279			
General	7286.928	398.351	9	.000

The null hypothesis states that the location parameters (slope coefficients) are the same across response categories.

a. Link function: Logit.

68