## LOGISTIC REGRESSION

Chapter 6

## Nominal Dependent Variables

There are three main approaches to handling nominal dependent variables:

- 1. Discriminant analysis (Chapter 5).
  - Dependent variable has two or more categories
- 2. Logistic regression (Chapter 6).
  - Dependent variable has two categories.
- 3. Multinomial logistic regression.
  - Dependent variable has more than two categories.

## Logistic Regression May Be Preferred . . .

Discriminant analysis:

Only continuous independent variables.

Finicky

- Needs equal variance—covariance matrices (Box's M test) across groups, and this assumptions is not met in many situations.
- Y can be nominal or ordinal.

Logistic regression:

- Continuous, ordinal, and nominal variables.
- Robust against assumption violations.
- Easier because it is similar to multiple regression, thus it is intuitively appealing.

### Example

We have a research question that requires a nominal dependent variable.

E.g. Foreign market entry mode of service firms can be divided between high control (subsidiary) and low control (agents etc.).

### **Converting Metric Variables to Non-Metric**

Did you start foreign operations because of demands from customers?

No 1 2 3 4 5 Yes

Polar extremes approach = compares only the extreme two groups and excludes the middle group.
 Median split = dividing the data at the median.

## Dichotomize a Variable



**Reviewer:** Why take a higher level variable and reduce it to a lower scale?

**Answer:** The variable was natural for dichotomization (I used polar extremes).

## What Logistic Regression Does

- Using maximum likelihood estimation, logistic regression predicts the probability of an event occurring, in this case foreign market entry mode.
- Contrary to "typical" multiple regression, it does not use ordinary least squares estimation.

## Logistic Curve



Predicts the probability of an event for levels of the independent variable(s).

## Logistic Curve



## Sample Size

- Maximum Likelihood Estimation requires larger samples than OLS.
- Hosmer Lemeshow recommend at least 400 plus hold-out sample.
- 10 observations per estimated parameter.
- Non-metric nominal or ordinal variables create more groups.

## Example "Choice of Foreign Market Entry Mode in Service Firms"

Anders Blomstermo D. Deo Sharma James Sallis

(2006), International Marketing Review, Vol. 23 No. 2, pp. 211-229.

## The Variables

Entry Mode: A dichotomous decision between a low control entry mode (coded 0), and a high control entry mode (coded 1).

Hard Soft: A dichotomous variable coded 0 for hard services and 1 for soft services.

**Relational friction:** It is easier and cheaper work in high control relationships (partner versus subsidiary).

**Experience:** A firm's previous experience in international markets.

Cultural distance: Hofstede (80) Kogut & Singh (88).

Firm size: Number of employees.

<ul> <li>Logistic Regression</li> <li>Hard 0 /Soft 1 service</li> <li>Relational friction [frict</li> <li>International experien</li> <li>Cultural distance [kul</li> <li>Firm size [size]</li> </ul>	Dependent: ✓ Entry mode [entry] Block 1 of 1 Previous Block 1 of 1 Next Block 1 of 1 hardsoft friction experinc undefined	X       Image: Categorical         Save       Statistics and Plots         Options       Classification plots         Style       Mosmer-Lemeshow goodness-of-fi         Bootstrap       Outliers outside       std. dev.         ØAll cases       Display         ØAt each step       At last step
	≥a*b> size <u>M</u> ethod: Enter ▼ Selection Varia <u>b</u> le:	Probability for Stepwise Entry: 0,05 Removal: 0,10
ОК	Paste Reset Cancel Help	Include con <u>s</u> tant in model

## Interpretation

Do model calculations based on variables in analysis.

Case Processing Summary						
Unweighted Case	s <sup>a</sup>	N	Percent			
Selected Cases	Included in Analysis	65	41,4			
	Missing Cases	92	58,6			

a. If weight is in effect, see classification table for the total

92

0

100,0

100,0

,0

157

157

Missing Cases

Total

Unselected Cases

number of cases.

Total

Dependent Variable Encoding

	Original Value	Internal Value			
	Low control	0			
	High control	1			
Dependent coding.					

## **Beginning Block**

#### Block 0: Beginning Block

	Entry mode				Percentage	
	Observed		Low control	High control	Correct	
Step 0	Entry mode	Low control	39	0	100,0	
		High control	26	0	0,	
	Overall Perce	entage			60,0	
			n = 65			

Classification Table<sup>a,b</sup>



## Model Fit

Good model fit is indicated by a significant model chi-square statistic and an insignificant Hosmer-Lemeshow chi-square statistic.

This is comparable to the F-statistic in OLS regression

		Chi-square	df	Sig.
Step 1	Step	22,206	5	,000,
	Block	22,206	5	,000,
	Model	22,206	5	,000,
				1

#### Omnibus Tests of Model Coefficients

# Good!

#### Hosmer and Lemeshow Test Step Chi-square df Sig. 1 11,800 7 ,107



## **Comparing Models**

## Refer to the -2 Log likelihood. The lower the number the better the fit.



## Hit Rate (classification accuracy)



## **Classification Accuracy**

The classification accuracy is measured by the hit-rate:

The naive model (maximum chance) is simply the accuracy if all observations were placed into the largest group. In the current example, the percentage correct would need to exceed 60%. And, the rule of thumb is that it should exceed by a margin of 25%.

60% x 1.25 = 75%

Since 87.2% > 75%, the naïve (maximum chance) model exceeds prediction by a good margin.

	Observed		Entry Low control	mode High control	Percentage Correct	
Step 1	Entry mode	Low control	34	5	87,2	60% Low control
		High control	8	18	69,2	40% High control
	Overall Perce	entage			80,0	2

#### Classification Table<sup>a</sup>

## **Classification Accuracy**

Another assessment of classification accuracy is the random model (proportional chance). It estimates how well the model classifies observations to the small group. It is calculated as:

 $\alpha^2$  + (1- $\alpha$ )<sup>2</sup>, where  $\alpha$  is the proportion in the small group.

In the current example, the small group proportion is 40%  $.4^2 + (1 - .4)^2 = .52$ . And, the rule of thumb is that it should exceed by a margin of 25%.

52% x 1.25 = 65%

Since 69.2 > 65, the random (proportional chance) model exceeds prediction by a good margin.

## Hypothesis Tests

The Wald statistic indicates the significance of each estimated coefficient, providing tests for individual hypotheses. See the Sig. (p-value)

• In "typical" multiple regression this is done with the t-tests.

		Variables in the Equation				$\wedge$	
		В	S.E.	Wald	df	Sig.	Exp(B)
Step 1 <sup>a</sup>	Hard 0 /Soft 1 services	1,725	,699	6,098	1	,014	5,612
	Relational friction	-,532	,340	2,440	1	,118	,588
	International experience	-,455	,235	3,746	1	,053	,634
	Cultural distance	,638	,246	6,739	1	,009	1,893
	Firm size	,000,	,000	,521	1	,471	1,000
	Constant	,907	1,654	,300	1	,584	2,476

 a. Variable(s) entered on step 1: Hard 0 /Soft 1 services, Relational friction, International experience, Cultural distance, Firm size.

## Coefficients

In logistic regression, the B coefficients measure the change in the odds of group membership on the dependent variable. They are difficult to directly interpret, so instead we look at the exponentiated logistic coefficient.

When the B is negative, the Exp(B) will be less than one. That means it reduces the odds of belonging to the dependent variable category coded 1.

When B is positive, the Exp(B) will be greater than one. That means it increases the odds of belonging to the dependent variable category coded 1.

## Magnitude of the Relationship ...

Metric variables, a one-unit change in the independent variable:

Exp(B) = 1, means no change in the odds.

Exp(B) = .5, means a 50% decrease in the odds of category 1.

Exp(B) = 2, means a 100% increase in the odds of category 1.

Dummy variables have the same interpretation as metric variables, but there is only the possibility of 0 or 1, not a one-unit change.

## Interpretation

		В	S.E.	Wald	df	Sig.	Exp(B)
Step 1 <sup>a</sup>	Hard 0 /Soft 1 services	1,725	,699	6,098	1	,014	5,612
	Relational friction	-,532	,340	2,440	1	,118	,588
	International experience	-,455	,235	3,746	1	,053	,634
	Cultural distance	,638	,246	6,739	1	,009	1,893
	Firm size	,000,	,000,	,521	1	,471	1,000
	Constant	,907	1,654	,300	1	,584	2,476

#### Variables in the Equation

 a. Variable(s) entered on step 1: Hard 0 /Soft 1 services, Relational friction, International experience, Cultural distance, Firm size.

Soft service firms (coded 1) are 5.612 times more likely to choose a high control entry mode (coded 1) than hard service firms.

A one-unit increase in cultural distance will increase the odds of high control entry by 89.3%.

A two-unit increase in cultural distance will increase the odds of high control entry by 178.6%.