

# Limited Dependent Variables

**Ciaran S. Phibbs**

**April 27, 2019**

# Objective

- This is a broad topic; one that doesn't get much attention in standard courses, except maybe logistic regression
  - This course will provide a brief overview of some other models
  - Intent is to introduce issues/models, not an in-depth treatment of them
-

# Limited Dependent Variables

- 0-1, small number of options, small counts, etc.
  - The dependent variable is not continuous, or even close to continuous.
-

# Types of Dependent Variables

- Binary Choice
  - Multinomial Choice (more than one option, can be ordered or unordered)
  - Counts
  - Most models in general framework of probability models
    - Prob (event/occurs)
-

# Binary Outcomes

## Common in Health Care

- Mortality
  - Other outcome
    - Infection
    - Patient safety event
    - Rehospitalization <30 days
  - Decision to seek medical care
-

$$Y_i = \beta_0 + \beta X + \varepsilon_i$$

$Y_i=0$  if lived,  $Y_i=1$  if died

$$\text{Prob}(Y_i=1) = F(X, \beta)$$

$$\text{Prob}(Y_i=0) = 1 - F(X, \beta)$$

OLS, also called a linear probability model

$\varepsilon_i$  is heteroscedastic, depends on  $\beta X$

Predictions not constrained to match actual outcome (0,1) and can get negative predicted values

---

# Standard Approaches to Binary Choice-1

- Logistic regression

$$\textit{prob}(Y=1) = \frac{e^{\beta X}}{1 + e^{\beta X}}$$

---

# Advantages of Logistic Regression

- Designed for relatively rare events
  - Commonly used in health care; most readers familiar with odds ratios
-

# Standard Approaches to Binary Choice-2

- Probit regression (classic example is decision to make a large purchase)

$$y^* = \beta X + \varepsilon$$

$$y=1 \text{ if } y^* > 0$$

$$y=0 \text{ if } y^* \leq 0$$

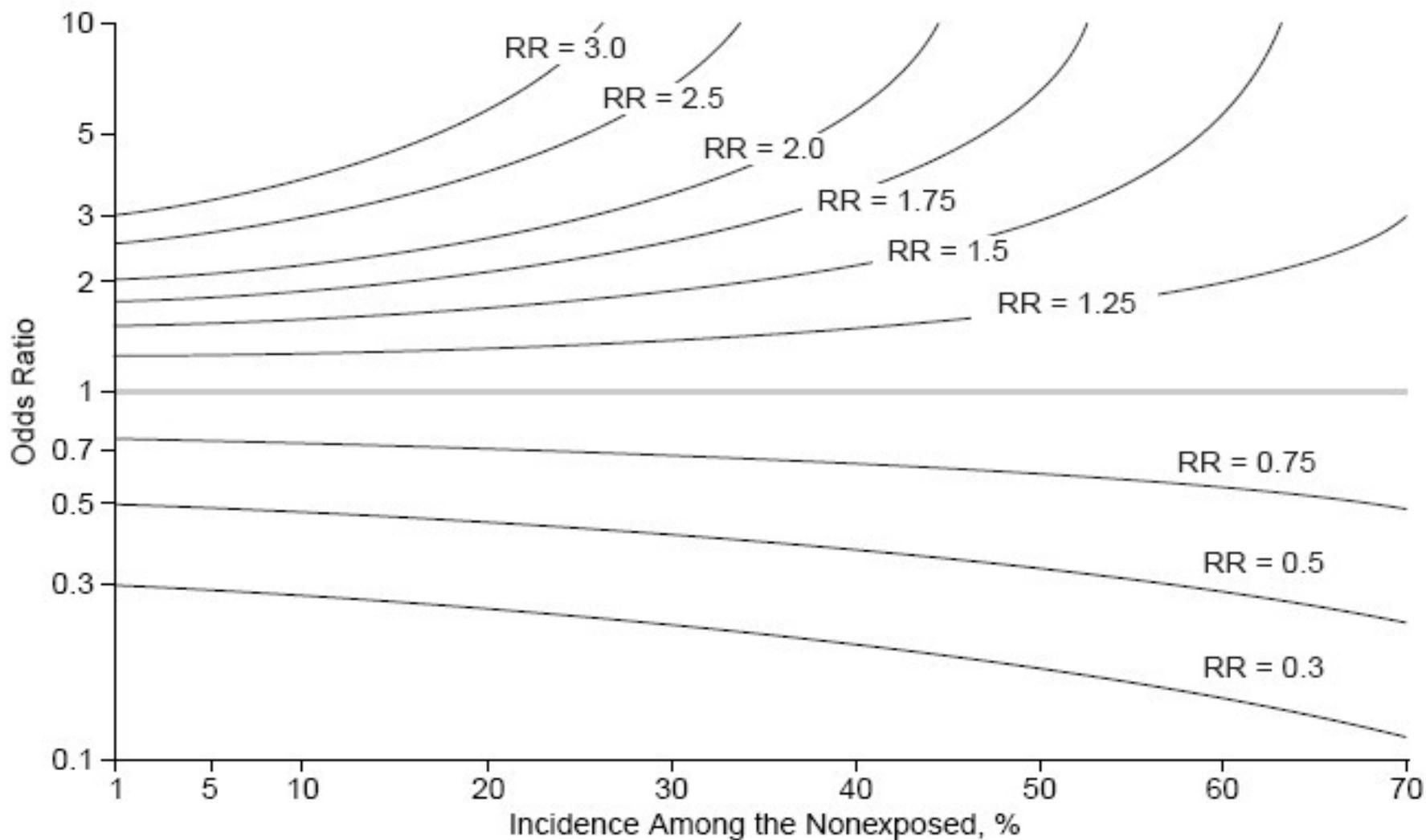
---

# Binary Choice

- There are other methods, using other distributions.
  - In general, logistic and probit give about the same answer.
  - It used to be a lot easier to calculate marginal effects with probit, not so any more
-

# Odds Ratios vs. Relative Risks

- Standard method of interpreting logistic regression is odds ratios.
  - Convert to % effect, really relative risk
  - This approximation starts to break down at 10% outcome incidence
-



The relationship between risk ratio (RR) and odds ratio by incidence of the outcome.

# Can Convert OR to RR

- Zhang J, Yu KF. What's the Relative Risk? A Method of Correcting the Odds Ratio in Cohort Studies of Common Outcomes. JAMA 1998;280(19):1690-1691.

$$\mathbf{RR = \frac{OR}{(1-P_0) + (P_0 \times OR)}}$$

Where  $P_0$  is the sample probability of the outcome

---

# Effect of Correction for RR

From Phibbs et al., NEJM 5/24/2007,  $\approx 20\%$  mortality

| <b>Odds Ratio</b> | <b>Calculated RR</b> |
|-------------------|----------------------|
| <b>2.72</b>       | <b>2.08</b>          |
| <b>2.39</b>       | <b>1.91</b>          |
| <b>1.78</b>       | <b>1.56</b>          |
| <b>1.51</b>       | <b>1.38</b>          |
| <b>1.08</b>       | <b>1.06</b>          |

---

# OR vs. RR

- Zhang is an approximation, not exact.
  - Many journals, especially epidemiology journals, now want direct estimation of RR
  - One option, Poisson with robust error variance
    - For binary outcome, IRR is a RR
    - Can run in most statistical packages; Stata has a Poisson command, or SAS w/ Proc GENMOD
-

# OR vs. RR

- Except for very large OR, and when incidence rate is large, effects of correct estimation of RR are relatively modest.
  - But, better to do it correctly, then know that you won't be over-estimating the effect.
-

# Extensions for Binary Data

- There are a lot of variations
    - Panel data
    - Grouped data
  - For panel data, can now estimate both random effects and fixed effects models.
  - The Stata manual lists over 30 related estimation commands for binary outcomes
-

# Extensions

- Goodness of fit tests. Several tests.
  - Probably the most commonly reported statistics are:
    - Area under ROC curve, c-statistic in SAS output. Range 0.50 to 1.0. Intuitively, how well does model predict compared to random assignment
    - Hosmer-Lemeshow test of goodness of fit
-

# More on Hosmer-Lemeshow Test

- The H-L test breaks the sample up into  $n$  (usually 10, some programs (Stata) let you vary this) equal groups and compares the number of observed and expected events (e.g., deaths) in each group.
  - If your model predicts well, the events will be concentrated in the highest risk groups; most can be in the highest risk group.
-

# More on Hosmer-Lemeshow Test

- There can be cases when you want to understand how well the model predicts in the tail, i.e., how well does the model discriminate in the mortality prediction, where most deaths are in the top decile because you have a good model?
  - Alternate specification, divide the sample so that the events are split into equal groups (deciles), instead of the number of cases, compare the observed and expected deaths in each decile. Not a valid formal test, but useful to learn more about how well you are actually predicting outcome in the tail.
-

# Estimation Note for Very Large Samples

- If you have very large samples; millions, it takes a lot longer to estimate a maximum likelihood model than OLS
  - But, same  $X$  matrix, so the p-values for OLS are approximately the same as a logit model. Can use OLS for model development, test RHS variables as discussed last week, and only estimate the final models with logit or other maximum likelihood model.
-

# Multinomial or Discrete Choice

- What if more than one choice or outcome?
  - Options are more limited
    - Multivariable Probit (multiple decisions, each with two alternatives)
    - Two different logit models (single decision, multiple alternatives)
    - Ordered logit
-

# **Examples of Health Care Uses for Multiple Choice Models**

- Choice of what hospital to use, among those in market area (or chose VA vs. several other options)
  - Choice of treatment among several options
  - Ordered vs. unordered choices
-

# Logit Models for Multiple Choices

- Conditional Logit Model (McFadden)
    - Unordered choices
  - Multinomial Logit Model
    - Choices can be ordered.
-

# Conditional logit model

- Also known as the random utility model
  - Is derived from consumer theory
  - How consumers choose from a set of options
  - Model driven by the characteristics of the choices; original McFadden model only included characteristics of choices, no individual characteristics.
  - Can express the results as odds ratios.
-

# Conditional logit model

- Individual characteristics “cancel out” but can be included indirectly. For example, in hospital choice, can interact individual characteristic with distance to hospital
  - Stata, clogit, can directly include both individual and choice-level variables.
-

# Estimation of McFadden's Model

- Some software packages (e.g. SAS) require that the number of choices be equal across all observations.
  - LIMDEP, allows a “NCHOICES” options that lets you set the number of choices for each observation. This is a very useful feature. Can to do this in Stata (clogit) with “group”
-

# Multinomial Logit Model

- Must identify a reference choice, model yields set of parameter estimates for each of the other choices, relative to the reference choice
  - Allows direct estimation of parameters for individual characteristics. Model can (should) also include parameters for choice characteristics
-

# Multinomial Logit Model

- Multinomial logit model can be estimated for either ordered or unordered choices. State (mlogit) allows up to 50 choices (or outcomes).
-

# Examples of Application of Multinomial Choice

- NICUs have formal levels of care that define patients that they can treat
  - Compare outcomes of infants born in lower level NICUs, compared to outcomes of infants born in highest level.
  - Full example in clinical journal, Haberland et al., *Pediatrics* 2006;118(6):1667-1675.
-

# Independence of Irrelevant Alternatives

- Results should be robust to varying the number of alternative choices
    - Can re-estimate model after deleting some of the choices.
    - McFadden, regression based test. Regression-Based Specification Tests for the Multinomial Logit Model. *J Econometrics* 1987;34(1/2):63-82.
  - If fail IIA, may need to estimate a nested logit model
-

# Independence of Irrelevant Alternatives - 2

- McFadden test can also be used to test for omitted variables.
  - For many health applications, doesn't matter, the models are very robust (e.g. hospital choice models driven by distance).
-

# Count Data (integers)

- Continuation of the same problem; dependent variable can only assume specific values and can't be  $< \text{zero}$
  - Problem diminishes as counts increase
  - Rule of Thumb. Need to use count data models for counts under 30
-

# Count Data

- Some examples of where count data models are needed in health care
    - Dependent variable is number of outpatient visits by each patient
    - Number of times a prescription for a chronic disease medication is refilled in a year
    - Number of adverse events in a unit (or hospital) over a period of time
-

# Count Data

- Poisson distribution. A distribution for counts.
  - Problem: very restrictive assumption that mean and variance are equal

$$\text{prob}(Y_i = y_i) = \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!}$$

---

# Count Data

- In general, negative binomial is a better choice. Stata (nbgreg), test for what distribution is part of the package. Other distributions can also be used.

$$f(y_i | x_i, u_i) = \frac{e^{-\lambda u_i} (\lambda u_i)^{y_i}}{y_i!}$$

---

# Interpreting Count Data Models

- $\ln E(\text{event rate}) = Bx$
  - Incidence Rate Ratio =  $e^B$ , similar to an odds ratio; rate at which event occurs
  - Multiply IIR by exposure to get expected number of events
-

# Notes for Count Data Models

- More common to see OLS used for counts than for binary or very limited choices.
- Real problem with OLS when there are lots of zeros. Will result in reduced statistical significance. Can go in opposite direction when counts are larger.
- Can have meaningful effects of coefficients

# Notes for Count Data Models-2

- As with logit, lots of extensions, for example, Stata, xtnbreg for panel data
  - Stata, gnbreg, generalization of negative binomial model with respect to dispersion of the data
-

# Notes for Count Data Models-3

- 30 is a rule of thumb, but should still consider a count model if most are small counts
  - Need to consider distribution and data generating process. If mixed process, may need to split sample
-

# Example of Mixed Data Generating Processes

- Predicting LOS for newborns
    - Well babies, all with  $LOS \leq 5$  days, about 90% of all newborns, clearly a count
    - Sick newborns, can have very long LOS
  - Solution, separate models for well babies and for sick babies.
-

# Other Models

- New models are being introduced all of the time. More and better ways to address the problems of limited dependent variables.
  - Includes semi-parametric and non-parameteric methods.
-

# Reference Texts

- Greene. Econometric Analysis
  - Wooldridge. Econometric Analysis of Cross Section and Panel Data
  - Maddala. Limited-Dependent and Qualitative Variables in Econometrics
-

# Journal References

- McFadden D. Specification Tests for the Multinomial Logit Model. *J Econometrics* 1987;34(1/2):63-82.
  - Zhang J, Yu KF. What's the Relative Risk? A Method of Correcting the Odds Ratio in Cohort Studies of Common Outcomes. *JAMA* 1998;280(19):1690-1691.
-

# Next lecture

Fixed Effects & Random Effects

Liam Rose

April 3, 2019