

Limited Dependent Variables & Selection: PS #1

Francis DiTraglia

HT 2021

This problem set is on *Friday in Week 2 of HT 2021*. You need only submit solutions to questions 1–3, as question #4 will not be marked. See the explanation immediately preceding question #4 for further information.

- Let $y \sim \text{Poisson}(\theta)$.
 - Using steps similar to the derivation of $\mathbb{E}[y]$ from the lecture slides, show that $\mathbb{E}[y(y-1)] = \theta^2$.
 - Use your answer to the preceding part, along with the result $\mathbb{E}[y] = \theta$, to show that $\text{Var}(y) = \theta$.
- Suppose that we observe count data $y_1, \dots, y_N \sim \text{iid } p_\theta$ and our model $f(y_i|\theta)$ is a $\text{Poisson}(\theta)$ probability mass function. Show that $\hat{K} = s_y^2/(\bar{y})^2$ where we define $s_y^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2$ and $\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$.
- Let $\hat{\beta}$ be the conditional maximum likelihood estimator of β_o in a Poisson regression model with conditional mean function $\mathbb{E}(y_i|\mathbf{x}_i) = \exp(\mathbf{x}_i'\beta_o)$, based on a sample of iid observations $(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)$.
 - Derive the first-order conditions for $\hat{\beta}$.
 - Using your answer to the previous part show that, so long as \mathbf{x}_i includes a constant, the residuals $\hat{u}_i \equiv y_i - \exp(\mathbf{x}_i'\hat{\beta})$ sum to zero, as in OLS regression.
 - Using your answer to the preceding part, show that $\left[\frac{1}{N} \sum_{i=1}^N \exp(\mathbf{x}_i'\hat{\beta}) \right] = \bar{y}$, where \bar{y} is the sample mean of y , so that $\bar{y}\hat{\beta}_j$ equals the estimated average partial effect of x_j in this model.
 - Explain why multiplying the estimated coefficients from this model by \bar{y} makes them roughly comparable to the corresponding OLS estimates from the model $y_i = \mathbf{x}_i'\theta + \varepsilon_i$.

The following applied question will *not be marked*, but you are encouraged to complete it nonetheless as it will build your understanding of the material from the lectures. Solving this problem will require some of the R material from Lecture #6.

- This question is adapted from Wooldridge (2010)*. To answer it you will need to use the dataset `SMOKE.RAW`, which can either be downloaded from the MIT Press website

for the text, or loaded directly into R using the package `Wooldridge`. Documentation for the dataset is available in the R package or alternatively at <http://fmwww.bc.edu/ec-p/data/wooldridge/smoke.des>

- (a) Use a linear regression to predict *cigs*, the number of cigarettes smoked each day, using the regressors $\log(\text{cigpric})$, $\log(\text{income})$, *restaurn*, *white*, *educ*, *age*, and age^2 . Interpret your findings. In particular: are cigarette prices and income statistically significant predictors? Does this depend on whether you use robust standard errors?
- (b) Repeat the preceding part but estimate a *Poisson* regression with an exponential conditional mean function rather than a linear regression. Calculate the APEs for the Poisson model and compare them to the OLS estimates.
- (c) If you calculated standard errors using the Poisson variance assumption, are cigarette prices and income statistically significant? Compare to your OLS results from above.
- (d) Calculate $\hat{\sigma}^2$. Does your estimate suggest evidence of overdispersion? If you use the Quasi-Poisson Variance assumption, how do your results compare to those of the preceding part?
- (e) How do your answers to the preceding two parts change if you instead use the fully-robust “sandwich” standard errors?