

# poisson: Some convergence issues

J. M. C. Santos Silva  
University of Essex and Centre for Applied Mathematics and Economics  
Colchester, United Kingdom  
jmcss@essex.ac.uk

Silvana Tenreyro  
London School of Economics Department of Economics  
Centre de Recerca en Economia Internacional, Center for Economic Performance,  
and Center for Economic and Policy Research  
London, United Kingdom  
s.tenreyro@lse.ac.uk

**Abstract.** In this article, we identify and illustrate some shortcomings of the `poisson` command in Stata. Specifically, we point out that the command fails to check for the existence of the estimates, and we show that it is very sensitive to numerical problems. While these are serious problems that may prevent users from obtaining estimates or may even produce spurious and misleading results, we show that the informed user often has simple workarounds available for addressing these problems.

**Keywords:** `st0225`, `ppml`, Poisson regression, collinearity, complete separation, numerical problems, perfect prediction

## 1 Introduction

Besides being the most widely used estimator for count data (see [Winkelmann \[2008\]](#) and [Cameron and Trivedi \[1998\]](#)), Poisson regression is also becoming increasingly used to estimate multiplicative models for other nonnegative data (see, among others, [Manning and Mullahy \[2001\]](#) and [Santos Silva and Tenreyro \[2006\]](#)). The availability in Stata of a command that estimates Poisson regression has been an important reason for the increasing popularity of this estimator.

However, researchers using Poisson regression, especially those using it to estimate gravity equations as recommended by [Santos Silva and Tenreyro \(2006\)](#), often find that the algorithm implemented in Stata's `poisson` command does not converge. There are two main reasons for this lack of convergence. First, as noted by [Santos Silva and Tenreyro \(2010\)](#), there are instances in which the estimates do not exist, and if that is the case, the convergence of the algorithm used to maximize the likelihood function can only be spurious. Second, even when the estimates exist, researchers using Stata may have trouble getting Poisson regression estimates because the `poisson` command is very sensitive to numerical problems.

In this article, we describe how researchers can identify some of the situations that may lead to convergence problems, and we propose some simple workarounds.

## 2 The nonexistence of the estimates

Let  $y_i$  and  $x_i$ , respectively, denote the variate of interest and the vector of covariates, and assume that the researcher specifies  $E(y_i|x_i) = \exp(x_i'\beta)$ . In a sample of size  $n$ ,  $\hat{\beta}$  (the Poisson regression estimate of  $\beta$ ) is defined by

$$\sum_{i=1}^n \left\{ y_i - \exp(x_i'\hat{\beta}) \right\} x_i = 0 \quad (1)$$

The form of (1) makes clear that  $\beta$  will be consistently estimated as long as the conditional mean is correctly specified. That is, the only condition required for the consistency of the estimator is that  $E(y_i|x_i) = \exp(x_i'\beta)$ . This is the well-known pseudo-maximum-likelihood result of [Gourieroux, Monfort, and Trognon \(1984\)](#).

However, [Santos Silva and Tenreyro \(2010\)](#) have shown that  $\hat{\beta}$  does not always exist and that its existence depends on the data configuration. In particular, the estimates may not exist if there is perfect collinearity for the subsample with positive observations of  $y_i$ .<sup>1</sup> If the estimates do not exist, either it is impossible for the estimation algorithm to converge or convergence is spurious. The following Stata code illustrates the situation where convergence is not achieved.<sup>2</sup>

```
drawnorm x1, n(1000) seed(101010) double clear
generate double u=rpoisson(1)
generate y=exp(1+10*x1)*u
generate double x2=(y==0)
poisson y x1 x2, vce(robust)
```

An example where the convergence is spurious is given by the code below:

```
drawnorm x1, n(1000) seed(101010) double clear
generate double y=rpoisson(1)
generate double x2=(y==0)
poisson y x1 x2, vce(robust)
```

The nonexistence of the maximum likelihood estimates in Poisson regression is analogous to what happens in binary choice models when there is complete separation or quasi-complete separation, as described by [Albert and Anderson \(1984\)](#) and [Santner and Duffy \(1986\)](#). In the case of binary models, it is standard to check for the existence of the estimates before starting the actual estimation. In contradistinction, the `poisson` command in Stata does not check for the existence of the estimates, and

1. See also [Haberman \(1973\)](#).

2. The code used in this article produces the desired results in Stata/IC 11.2 for Windows (32-bit). Using other flavors of Stata—for example, Stata/MP—may lead to different outcomes.

therefore it is important that users investigate whether the estimates exist. Because the regressors that may cause the nonexistence of the estimates are characterized by their perfect collinearity with the others for the subsample with  $y_i > 0$ , they can easily be identified in Stata by using a simplified, two-step version of the three-step method suggested by Santos Silva and Tenreyro (2010):

Step 1: Construct a subset of explanatory variables—say,  $\tilde{x}_i$ —comprising only the regressors that are not collinear for the observations with  $y_i > 0$ .

Step 2: Using the full sample, run the Poisson regression of  $y_i$  on  $\tilde{x}_i$ .

The following code, which assumes that all variables with names starting with **x** are regressors, illustrates the implementation of the procedure:<sup>3</sup>

```
local _rhs "x*"
_rmcoll `'_rhs'` if y>0
poisson y `r(varlist)`, vce(robust)
```

This procedure ensures that the estimates exist by eliminating all potentially problematic regressors, even those that actually do not lead to the nonexistence of the maximum likelihood estimates.<sup>4</sup> Therefore, the researcher should subsequently investigate one by one all the variables that were dropped to see if any of them can be included in the model. Careful investigation of the variables to be excluded is particularly important when the model contains sets of dummies with several categories; in this case, dropping one of the dummies implies an arbitrary redefinition of the reference category, which is unlikely to be sensible. In any case, dropping some regressors should never be an automatic procedure because it changes the model specification; therefore, the researcher should carefully consider what is the best way to find an interesting specification for which the (pseudo) maximum likelihood estimates exist.

The nonexistence of the estimates can also occur in any regression model where the conditional mean is specified in such a way that its image does not include all the points in the support of the dependent variable. Therefore, unless the data are strictly positive, this problem can occur not only in the Poisson regression but also in other models specifying  $E(y_i|x_i) = \exp(x_i'\beta)$ , and in models for limited dependent variables like the tobit (Tobin 1958). In all these cases, the identification of the problematic regressors can be done using methods akin to the one described above.

### 3 Numerical difficulties

Even if the (pseudo) maximum likelihood estimates of the Poisson regression exist, Stata may have difficulty identifying them because of the sensitivity to numerical problems of the algorithms available in the `poisson` command. In particular, we are aware of three situations in which the algorithms in the `poisson` command have trouble locating the

3. We are grateful to Markus Baldauf for help with the development of an earlier version of this code and to an anonymous referee for suggesting this much simpler version using the `_rmcoll` command.

4. A less strict criterion for selecting the regressors to be dropped is used by default in the `ppml` command, which is briefly discussed in section 4.

maximum and may not converge, even when the (pseudo) maximum likelihood estimates of the Poisson regression are well defined.

The simplest case in which Stata finds it difficult to find the Poisson (pseudo) maximum-likelihood estimates is when  $y$  has some very large values. The following Stata code illustrates the situation:<sup>5</sup>

```
drawnorm u x1 x2, n(1000) seed(101010) double clear
generate double y = exp(35+x1+x2+u)
poisson y x1 x2, vce(robust) difficult
```

In this example, the Poisson regression does not converge, at least not in a reasonable number of iterations. Obviously, in this case the problem can easily be bypassed just by rescaling the dependent variable, say, by dividing it by  $\exp(35)$ .

A second situation in which Stata finds it difficult to locate the solution of (1) occurs when the regressors are highly collinear and have very different magnitudes. The following Stata code illustrates the situation:<sup>6</sup>

```
drawnorm u e x1, n(1000) seed(101010) double clear
generate x2 = (x1<-2)
generate double x3 = 20+x1+(e/100)*(x1<-2)
generate double y = exp(1+x1+x2+u)
poisson y x1 x2 x3, vce(robust) difficult
```

In this case, again, the Poisson regression does not converge but a simple workaround is available: if the third regressor is recentered at zero, convergence is achieved with ease.

These two examples suggest that when facing convergence problems, researchers should rescale and recenter their data in a way that reduces possible numerical problems. However, even if that is done, Stata will have trouble finding the (pseudo) maximum likelihood estimates of the Poisson regression when the covariates are extremely (but not perfectly) collinear. The following example illustrates this situation:

```
drawnorm u e x1, n(1000) seed(101010) double clear
generate double x2 = (x1+e/18000)
generate double y = exp(1+x1+x2+u)
poisson y x1 x2, vce(robust) difficult
```

In cases like this, it is generally not possible to bypass the problem using some sort of data transformation, and different workarounds are needed.<sup>7</sup>

---

5. We are grateful to Alexandros Theloudis for showing us a dataset where this situation occurs.

6. We are grateful to Avni Hanedar for showing us a dataset where this situation occurs.

7. Of course, the researcher may want to reconsider the specification being used.

## 4 Workaround

When the (pseudo) maximum likelihood estimates exist but convergence is not achieved with the default options, an obvious alternative to explore is to try one of the different optimization methods offered by the `poisson` command. However, for instance in the third example in section 3, none of the methods available leads to satisfactory results. Indeed, in that case, with the `technique(nr)` and the `technique(bhhh)` options, the algorithm fails to converge; and with the `technique(dfp)` and the `technique(bfgs)` options, the algorithm converges to a result that is far from the optimum. Alternatively, one can ensure convergence just by relaxing the convergence criteria. This, however, is a risky option because the algorithm may be stopped too soon, therefore not delivering the desired (pseudo) maximum likelihood estimates. This is what happens, for instance, when the `nonrtol` option is used.

A simple workaround that often (but by no means always) works is to use the `glm` command with the options `family(poisson)`, `link(log)`, and `irls`. Indeed, the iterated, reweighted least-squares algorithm provided by the `glm` command appears to be much more stable than the algorithms available in the `poisson` command, and it produces the correct results in the three examples presented in section 3.

To facilitate the estimation of Poisson regressions while Stata does not improve the reliability of `poisson`, we have written the `ppml` command, which checks for the existence of the (pseudo) maximum likelihood estimates and offers two methods to drop regressors that may cause the nonexistence of the estimates. Estimation is then implemented using the `glm` method, and `ppml` warns if the variables have large values that are likely to create numerical problems or if there are signs that the convergence is spurious.<sup>8</sup> Further details on `ppml` can be found in the corresponding help file.

## 5 Conclusions

In this article, we illustrated some shortcomings of Stata's `poisson` command. We believe that it should be relatively easy to update the `poisson` command so that it checks for the existence of the Poisson regression estimates and is more resilient to numerical problems.

Although an upgraded version of `poisson` is not available, practitioners can use our `ppml` command, which checks for the existence of the estimates before trying to estimate a Poisson regression and provides several warnings about possible convergence problems.

---

8. A telltale sign that the convergence is spurious is that some zero observations of  $y$  are “perfectly predicted”; in the second example in section 2, the values of  $\exp(x_i'\hat{\beta})$  for  $y = 0$  vary between  $3.83\text{e-}09$  and  $3.85\text{e-}09$ .

## 6 Acknowledgments

We are very grateful to an anonymous referee for most-valuable comments and suggestions. We also thank Styliani Christodouloupoulou for helpful discussions and comments on an earlier version of the paper. The usual disclaimer applies. Santos Silva gratefully acknowledges partial financial support from Fundação para a Ciência e a Tecnologia (FEDER/POCI 2010).

## 7 References

- Albert, A., and J. A. Anderson. 1984. On the existence of maximum likelihood estimates in logistic models. *Biometrika* 71: 1–10.
- Cameron, A. C., and P. K. Trivedi. 1998. *Regression Analysis of Count Data*. Cambridge: Cambridge University Press.
- Gourieroux, C., A. Monfort, and A. Trognon. 1984. Pseudo maximum likelihood methods: Applications to Poisson models. *Econometrica* 52: 701–720.
- Haberman, S. J. 1973. Log-linear models for frequency data: Sufficient statistics and likelihood equations. *Annals of Statistics* 1: 617–632.
- Manning, W. G., and J. Mullahy. 2001. Estimating log models: to transform or not to transform? *Journal of Health Economics* 20: 461–494.
- Santner, T. J., and D. E. Duffy. 1986. A note on A. Albert and J. A. Anderson's conditions for the existence of maximum likelihood estimates in logistic regression models. *Biometrika* 73: 755–758.
- Santos Silva, J. M. C., and S. Tenreyro. 2006. The log of gravity. *Review of Economics and Statistics* 88: 641–658.
- . 2010. On the existence of the maximum likelihood estimates in Poisson regression. *Economics Letters* 107: 310–312.
- Tobin, J. 1958. Estimation of relationships for limited dependent variables. *Econometrica* 26: 24–36.
- Winkelmann, R. 2008. *Econometric Analysis of Count Data*. 5th ed. Berlin: Springer.

### About the authors

João Santos Silva is an econometrician and professor of economics at the University of Essex. Silvana Tenreyro is a macroeconomist and reader in economics at the London School of Economics. The authors have done joint work on econometric models for nonnegative variables with a mass-point at zero, particularly on the estimation of the gravity equation for trade flows.