

# A Replication (and Tribute) of The Log of Gravity

Mauricio “Pachá” Vargas Sepúlveda

v. 2022-11-23 22:33

This work is licensed under a [Creative Commons “Attribution-NonCommercial-ShareAlike 4.0 International”](#) license.



## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Goal . . . . .	1
1.2	Usual disclaimer . . . . .	1
<b>2</b>	<b>Obtaining the original codes and data</b>	<b>1</b>
<b>3</b>	<b>Loading the original data</b>	<b>2</b>
<b>4</b>	<b>Replication attempt</b>	<b>2</b>
4.1	Poisson Pseudo Maximum Likelihood . . . . .	2
4.2	Ordinary Least Squares . . . . .	5
4.3	Tobit . . . . .	5
4.4	Non-Linear Least Squares . . . . .	7
<b>5</b>	<b>Replication results</b>	<b>7</b>
<b>6</b>	<b>References</b>	<b>9</b>

# 1 Introduction

## 1.1 Goal

Silva and Tenreyro<sup>1</sup> wasn't just an influential article, it defined my interest in gravity models to the point I wrote a master's thesis on it for UN ESCAP. Here I replicate the main results from the original article in R. The original results were obtained in Stata back in 2006.

The idea here is to be explicit regarding the conceptual approach to regression in R. For most of the replication I used base R<sup>2</sup> without external libraries (i.e packages) except when it was absolutely necessary.

Much of the methods exposed here lead to the exact same results as using the gravity package<sup>3</sup> which provides convenient wrappers for gravity estimation.

All questions are welcome to the address `mv.sepulveda@mail.utoronto.ca`.

## 1.2 Usual disclaimer

The views and opinions expressed in this course are solely those of the author and do not necessarily reflect the official position of any unit of the United Nations, the University of Toronto or the Pontifical Catholic University of Chile.

# 2 Obtaining the original codes and data

I shall organize the original codes and data from the authors' site to put these on GitHub and therefore ease reproducibility in case of broken links or anything that makes it difficult to obtain the original zip file with the data and codes.

```
url <- "https://personal.lse.ac.uk/tenreyro/regressors.zip"
zip <- gsub(".*/", "", url)
if (!file.exists(zip)) try(download.file(url, zip))
```

---

<sup>1</sup>“The Log of Gravity.”

<sup>2</sup>R Core Team, *R*.

<sup>3</sup>Woelwer et al., *Gravity*.

```
dout <- "regressors"
if (!dir.exists(dout)) unzip(zip, exdir = dout)
```

### 3 Loading the original data

Thanks to the `haven` package<sup>4</sup> we can read Stata datasets directly in R without loss of information about column types and other common problems when reading proprietary formats.

```
log_of_gravity <- haven::read_dta(paste0(dout, "/Log of Gravity.dta"))
```

## 4 Replication attempt

### 4.1 Poisson Pseudo Maximum Likelihood

Table 3 in Silva and Tenreyro<sup>5</sup> summarises a large portion of the article and it can be partially replicated with the following Stata code for the Poisson Pseudo Maximum Likelihood.

```
ppml trade lypex lypim lyex lyim ldist border comlang colony landl_ex
      landl_im lremot_ex lremot_im comfrt_wto open_wto
```

In R we would replicate it by fitting two Generalized Linear Models since the article introduces estimates with and without removing zero flows.

```
ppml_formula <- trade ~ lypex + lypim + lyex + lyim + ldist + border +
  comlang + colony + landl_ex + landl_im + lremot_ex + lremot_im +
  comfrt_wto + open_wto
```

---

<sup>4</sup>Wickham and Miller, *haven: Import and Export SPSS, Stata and SAS Files*.

<sup>5</sup>“The Log of Gravity.”

```

fit_ppml_1 <- glm(
  ppml_formula,
  data = log_of_gravity,
  subset = trade > 0,
  family = quasipoisson()
)

fit_ppml_2 <- glm(
  ppml_formula,
  data = log_of_gravity,
  family = quasipoisson()
)

kable(broom::tidy(fit_ppml_1))

```

term	estimate	std.error	statistic	p.value
(Intercept)	-31.5295516	0.5962458	-52.8801207	0.0000000
lypex	0.7213276	0.0081581	88.4190011	0.0000000
lypim	0.7318762	0.0081843	89.4248659	0.0000000
lyex	0.1544318	0.0130506	11.8333558	0.0000000
lyim	0.1326857	0.0129337	10.2589262	0.0000000
ldist	-0.7763158	0.0177657	-43.6974358	0.0000000
border	0.2023691	0.0340575	5.9419810	0.0000000
comlang	0.7512811	0.0371767	20.2083736	0.0000000
colony	0.0199721	0.0432244	0.4620549	0.6440524
landl_ex	-0.8724183	0.0572665	-15.2343477	0.0000000
landl_im	-0.7034613	0.0540523	-13.0144515	0.0000000
lremot_ex	0.6471603	0.0476017	13.5953189	0.0000000
lremot_im	0.5492566	0.0475629	11.5480021	0.0000000

term	estimate	std.error	statistic	p.value
comfrt_wto	0.1794431	0.0356532	5.0330200	0.0000005
open_wto	-0.1394213	0.0385109	-3.6203123	0.0002958

```
kable(broom::tidy(fit_ppml_2))
```

term	estimate	std.error	statistic	p.value
(Intercept)	-32.3261029	0.4439044	-72.822223	0.0000000
lypex	0.7324808	0.0060759	120.554228	0.0000000
lypim	0.7410780	0.0061147	121.196857	0.0000000
lyex	0.1567117	0.0098043	15.983898	0.0000000
lyim	0.1350185	0.0097182	13.893327	0.0000000
ldist	-0.7838006	0.0132914	-58.970409	0.0000000
border	0.1929108	0.0255010	7.564828	0.0000000
comlang	0.7459840	0.0278186	26.816062	0.0000000
colony	0.0250065	0.0323900	0.772044	0.4400983
landl_ex	-0.8634737	0.0428137	-20.168181	0.0000000
landl_im	-0.6964204	0.0404198	-17.229694	0.0000000
lremot_ex	0.6598400	0.0356853	18.490539	0.0000000
lremot_im	0.5615002	0.0355874	15.778062	0.0000000
comfrt_wto	0.1811072	0.0266624	6.792610	0.0000000
open_wto	-0.1068187	0.0287357	-3.717287	0.0002020

The replication effort here is null, it just sufficed to look at the summary table in the article and subset the data to drop zero flows. Therefore, it makes sense to proceed with the other models.

## 4.2 Ordinary Least Squares

The only consideration here is to drop zero flows for some of the models with log in the dependent variable even when Table 3 is not explicit about this, otherwise we break the fitting algorithm.

For example, for estimations of the type  $\log(\text{trade}) = \beta_0 + \beta_1 \text{lypex} + \dots + \varepsilon$ , we need to drop zero flows to replicate the result. On the other hand, for estimations of the type  $\log(1 + \text{trade}) = \beta_0 + \beta_1 \text{lypex} + \dots + \varepsilon$ , we don’t need to drop zero flows.

```
fit_ols_1 <- lm(
  update.formula(ppml_formula, log(.) ~ .),
  data = log_of_gravity,
  subset = trade > 0
)

fit_ols_2 <- lm(
  update.formula(ppml_formula, log(1 + .) ~ .),
  data = log_of_gravity
)
```

## 4.3 Tobit

The Tobit estimation is similar but requires the use of the censReg package.<sup>6</sup> The complicated part of the estimation here is to extract the right hand side of the model formula to define a vector of zeroes of the length of this right hand side plus two as starting point for the Maximum Likelihood estimation (i.e including the depending variable and intercept besides the estimating slopes).

In order to obtain the  $a$  value that matches the results in the article I proceeded with an iteration loop until achieving convergence with respect to one of the estimated slopes. The initial value of  $a = 200$  was arbitrary and set after trying reasonable guesses that converge to

---

<sup>6</sup>Henningsen, *censReg: Censored Regression (Tobit) Models*.

the slopes in the original article after 9 iterations for a final value of  $a = 159$ .

```
a <- 200
lypex_ref <- 1.058
tol <- 0.001
lypex_estimate <- 2 * lypex_ref
iter <- 0

while (abs(lypex_estimate - lypex_ref) > tol) {
  log_of_gravity$log_trade_cens <- log(a + log_of_gravity$trade)
  log_trade_cens_min <- min(log_of_gravity$log_trade_cens, na.rm = TRUE)

  fit_tobit <- censReg::censReg(
    formula = update.formula(ppml_formula, log_trade_cens ~ .),
    left = log_trade_cens_min,
    right = Inf,
    data = log_of_gravity,
    start = rep(0, 2 + length(attr(terms(ppml_formula), "term.labels"))),
    method = "BHHH"
  )

  lypex_estimate <- coef(fit_tobit)[2]
  if (abs(lypex_estimate - lypex_ref) > 2 * tol) {
    a <- a - 5
  } else {
    a <- a - 1
  }
  iter <- iter + 1
}
```

## 4.4 Non-Linear Least Squares

For this type of estimation the starting values are retrieved from the results of the PPML model with zero flows and then we pass these values to a Generalized Linear Model using the Gaussian distribution and a log-link.

```
fit_ppml_eta <- fit_ppml_2$linear.predictors
fit_ppml_mu <- fit_ppml_2$fitted.values
fit_ppml_start <- fit_ppml_2$coefficients

fit_nls <- glm(
  ppml_formula,
  data = log_of_gravity,
  family = gaussian(link = "log"),
  etastart = fit_ppml_eta,
  mustart = fit_ppml_mu,
  start = fit_ppml_start,
  control = list(maxit = 200, trace = FALSE)
)
```

## 5 Replication results

There wasn't much effort involved in the replication, which is something desirable. I didn't even have to email the authors with questions whereas the data was filtered or transformed in ways not mentioned in the article, which is something that we often see. Unlike many articles, this is one of the very few articles that I've found which passes the reproducibility review and is very close to full replication according to the criteria from Peng<sup>7</sup>.

---

<sup>7</sup>“Reproducible Research in Computational Science.”



Table 3: Replication results for OLS (1-2), Tobit (3), NLS (4) and PPML (5-6).

	<i>Dependent variable:</i>					
	(1)	(2)	(3)	(4)	(5)	(6)
lypex	0.938*** (0.012)	1.128*** (0.011)	1.059*** (0.011)	0.738*** (0.004)	0.721*** (0.008)	0.732*** (0.006)
lypim	0.798*** (0.011)	0.866*** (0.011)	0.848*** (0.010)	0.862*** (0.005)	0.732*** (0.008)	0.741*** (0.006)
lyex	0.207*** (0.017)	0.277*** (0.017)	0.228*** (0.014)	0.396*** (0.010)	0.154*** (0.013)	0.157*** (0.010)
lyim	0.106*** (0.017)	0.217*** (0.017)	0.178*** (0.014)	−0.033*** (0.007)	0.133*** (0.013)	0.135*** (0.010)
ldist	−1.166*** (0.034)	−1.151*** (0.037)	−1.160*** (0.029)	−0.924*** (0.008)	−0.776*** (0.018)	−0.784*** (0.013)
border	0.314** (0.143)	−0.241 (0.164)	−0.225** (0.109)	−0.081*** (0.010)	0.202*** (0.034)	0.193*** (0.026)
comlang	0.678*** (0.064)	0.742*** (0.064)	0.759*** (0.052)	0.689*** (0.016)	0.751*** (0.037)	0.746*** (0.028)
colony	0.397*** (0.068)	0.392*** (0.068)	0.416*** (0.056)	0.036** (0.018)	0.020 (0.043)	0.025 (0.032)
landl_ex	−0.062 (0.065)	0.106* (0.060)	−0.038 (0.060)	−1.367*** (0.031)	−0.872*** (0.057)	−0.863*** (0.043)
landl_im	−0.665*** (0.063)	−0.278*** (0.060)	−0.478*** (0.059)	−0.471*** (0.022)	−0.703*** (0.054)	−0.696*** (0.040)
lremot_ex	0.467*** (0.078)	0.526*** (0.089)	0.563*** (0.077)	1.188*** (0.018)	0.647*** (0.048)	0.660*** (0.036)
lremot_im	−0.205** (0.081)	−0.109 (0.089)	−0.032 (0.074)	1.010*** (0.018)	0.549*** (0.048)	0.562*** (0.036)
comfrt_wto	0.491*** (0.105)	1.289*** (0.143)	0.728*** (0.113)	0.443*** (0.014)	0.179*** (0.036)	0.181*** (0.027)
open_wto	−0.170*** (0.049)	0.739*** (0.048)	0.310*** (0.040)	0.928*** (0.024)	−0.139*** (0.039)	−0.107*** (0.029)
logSigma			0.677*** (0.007)			
Constant	−28.492*** (1.088)	−39.909*** (1.221)	−36.626*** (1.059)	−45.098*** (0.239)	−31.530*** (0.596)	−32.326*** (0.444)
Observations	9,613	18,360	18,360	18,360	9,613	18,360

Note:

v. 2022-11-23 22:33

\*p&lt;0.1; \*\*p&lt;0.05; \*\*\*p&lt;0.01

## 6 References

- Henningsen, Arne. *censReg: Censored Regression (Tobit) Models*, 2020. <https://CRAN.R-project.org/package=censReg>.
- Peng, Roger D. “Reproducible Research in Computational Science.” *Science* 334, no. 6060 (2011): 1226–27.
- R Core Team. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing, 2021. <https://www.R-project.org/>.
- Silva, JMC Santos, and Silvana Tenreyro. “The Log of Gravity.” *The Review of Economics and Statistics* 88, no. 4 (2006): 641–58.
- Wickham, Hadley, and Evan Miller. *haven: Import and Export SPSS, Stata and SAS Files*, 2021. <https://CRAN.R-project.org/package=haven>.
- Woelwer, Anna-Lena, Jan Pablo Burgard, Joshua Kunst, and Mauricio Vargas. *Gravity: Estimation Methods for Gravity Models*, 2020. <http://pacha.dev/gravity>.