poisson — Poisson regression

Description	Quick start
Options	Remarks and examples
References	Also see

Menu Stored results Syntax Methods and formulas

Description

poisson fits a Poisson regression of *depvar* on *indepvars*, where *depvar* is a nonnegative count variable.

If you have panel data, see [XT] **xtpoisson**.

Quick start

Poisson regression of y on x

poisson y x

Add categorical variable a

poisson y x i.a

Add exposure variable v

poisson y x i.a, exposure(v)

With robust standard errors

poisson y x i.a, vce(robust)

Report results as incidence-rate ratios

poisson y x i.a, irr

Replace data in memory with the results of running a Poisson regression model on each level of catvar statsby, by(catvar) clear: poisson y x

Menu

Statistics > Count outcomes > Poisson regression

Syntax

poisson depvar [indepvars] [if] [in] [weight] [, options]

options	Description
Model	
$\frac{nocons}{exposure(varname_e)}$ $\frac{off}{set(varname_o)}$ $\frac{const}{const}$	suppress constant term include $ln(varname_e)$ in model with coefficient constrained to 1 include $varname_o$ in model with coefficient constrained to 1 apply specified linear constraints
SE/Robust	
vce(<i>vcetype</i>)	<pre>vcetype may be oim, robust, cluster clustvar, opg, bootstrap, or jackknife</pre>
Reporting	
<u>l</u> evel(#)	set confidence level; default is level(95)
<u>ir</u> r	report incidence-rate ratios
<u>nocnsr</u> eport	do not display constraints
display_options	control columns and column formats, row spacing, line width, display of omitted variables and base and empty cells, and factor-variable labeling
Maximization	
maximize_options	control the maximization process; seldom used
<u>col</u> linear coeflegend	keep collinear variables display legend instead of statistics

indepvars may contain factor variables; see [U] 11.4.3 Factor variables.

depvar, indepvars, varname, and varname, may contain time-series operators; see [U] 11.4.4 Time-series varlists.

bayes, bayesboot, bootstrap, by, fmm, fp, jackknife, mfp, mi estimate, nestreg, rolling, statsby, stepwise, and svy are allowed; see [U] **11.1.10 Prefix commands**. For more details, see [BAYES] **bayes: poisson** and [FMM] **fmm: poisson**.

vce(bootstrap) and vce(jackknife) are not allowed with the mi estimate prefix; see [MI] mi estimate.

Weights are not allowed with the bootstrap prefix; see [R] bootstrap.

vce() and weights are not allowed with the svy prefix; see [SVY] svy.

fweights, iweights, and pweights are allowed; see [U] 11.1.6 weight.

collinear and coeflegend do not appear in the dialog box.

See [U] 20 Estimation and postestimation commands for more capabilities of estimation commands.

Options

[Model]

noconstant, exposure($varname_e$), offset($varname_o$), constraints(*constraints*); see [R] Estimation options.

SE/Robust

vce(vcetype) specifies the type of standard error reported, which includes types that are derived from asymptotic theory (oim, opg), that are robust to some kinds of misspecification (robust), that allow for intragroup correlation (cluster *clustvar*), and that use bootstrap or jackknife methods (bootstrap, jackknife); see [R] vce_option. Reporting

level(#); see [R] Estimation options.

irr reports estimated coefficients transformed to incidence-rate ratios, that is, e^{β_i} rather than β_i . Standard errors and confidence intervals are similarly transformed. This option affects how results are displayed, not how they are estimated or stored. irr may be specified at estimation or when replaying previously estimated results.

nocnsreport; see [R] Estimation options.

```
display_options: noci, nopvalues, noomitted, vsquish, noemptycells, baselevels,
allbaselevels, nofvlabel, fvwrap(#), fvwrapon(style), cformat(%fmt), pformat(%fmt),
sformat(%fmt), and nolstretch; see [R] Estimation options.
```

Maximization

```
maximize_options: difficult, technique(algorithm_spec), iterate(#), [no]log, trace,
gradient, showstep, hessian, showtolerance, tolerance(#), ltolerance(#),
nrtolerance(#), nonrtolerance, and from(init_specs); see [R] Maximize. These options are
seldom used.
```

Setting the optimization type to technique(bhhh) resets the default vcetype to vce(opg).

The following options are available with poisson but are not shown in the dialog box:

collinear, coeflegend; see [R] Estimation options.

Remarks and examples

The basic idea of Poisson regression was outlined by Coleman (1964, 378–379). See Cameron and Trivedi (2013; 2022, chap. 20) and Johnson, Kemp, and Kotz (2005, chap. 4) for information about the Poisson distribution. See Cameron and Trivedi (2013), Long (1997, chap. 8), Long and Freese (2014, chap. 9), McNeil (1996, chap. 6), and Selvin (2011, chap. 6) for an introduction to Poisson regression. Also see Selvin (2004, chap. 5) for a discussion of the analysis of spatial distributions, which includes a discussion of the Poisson distribution. An early example of Poisson regression was Cochran (1940).

Poisson regression fits models of the number of occurrences (counts) of an event. The Poisson distribution has been applied to diverse events, such as the number of soldiers kicked to death by horses in the Prussian army (von Bortkiewicz 1898); the pattern of hits by buzz bombs launched against London during World War II (Clarke 1946); telephone connections to a wrong number (Thorndike 1926); and disease incidence, typically with respect to time, but occasionally with respect to space. The basic assumptions are as follows:

- 1. There is a quantity called the *incidence rate* that is the rate at which events occur. Examples are 5 per second, 20 per 1,000 person-years, 17 per square meter, and 38 per cubic centimeter.
- 2. The incidence rate can be multiplied by exposure to obtain the expected number of observed events. For example, a rate of 5 per second multiplied by 30 seconds means that 150 events are expected; a rate of 20 per 1,000 person-years multiplied by 2,000 person-years means that 40 events are expected; and so on.
- 3. Over very small exposures ϵ , the probability of finding more than one event is small compared with ϵ .
- 4. Nonoverlapping exposures are mutually independent.

With these assumptions, to find the probability of k events in an exposure of size E, you divide E into n subintervals E_1, E_2, \ldots, E_n , and approximate the answer as the binomial probability of observing k successes in n trials. If you let $n \to \infty$, you obtain the Poisson distribution.

In the Poisson regression model, the incidence rate for the *j*th observation is assumed to be given by

$$r_i = e^{\beta_0 + \beta_1 x_{1,j} + \dots + \beta_k x_{k,j}}$$

If E_i is the exposure, the expected number of events, C_i , will be

$$C_j = E_j e^{\beta_0 + \beta_1 x_{1,j} + \dots + \beta_k x_{k,j}}$$
$$= e^{\ln(E_j) + \beta_0 + \beta_1 x_{1,j} + \dots + \beta_k x_{k,j}}$$

This model is fit by poisson. Without the exposure() or offset() options, E_j is assumed to be 1 (equivalent to assuming that exposure is unknown), and controlling for exposure, if necessary, is your responsibility.

Comparing rates is most easily done by calculating *incidence-rate ratios* (IRRs). For instance, what is the relative incidence rate of chromosome interchanges in cells as the intensity of radiation increases; the relative incidence rate of telephone connections to a wrong number as load increases; or the relative incidence rate of deaths due to cancer for females relative to males? That is, you want to hold all the x's in the model constant except one, say, the *i*th. The IRR for a one-unit change in x_i is

$$\frac{e^{\ln(E)+\beta_1x_1+\dots+\beta_i(x_i+1)+\dots+\beta_kx_k}}{e^{\ln(E)+\beta_1x_1+\dots+\beta_ix_i+\dots+\beta_kx_k}} = e^{\beta}$$

More generally, the IRR for a Δx_i change in x_i is $e^{\beta_i \Delta x_i}$. The lincom command can be used after poisson to display incidence-rate ratios for any group relative to another; see [R] lincom.

Example 1

Chatterjee and Hadi (2012, 174) give the number of injury incidents and the proportion of flights for each airline out of the total number of flights from New York for nine major US airlines in one year:

```
. use https://www.stata-press.com/data/r19/airline
```

. list

	airline	injuries	n	XYZowned
1.	1	11	0.0950	1
3.	3	7	0.0750	0
4.	4	19	0.2078	0
5.	5	9	0.1382	0
6.	6	4	0.0540	1
7.	7	3	0.1292	0
8.	8	1	0.0503	0
9.	9	3	0.0629	1

4

To their data, we have added a fictional variable, XYZowned. We will imagine that an accusation is made that the airlines owned by XYZ Company have a higher injury rate.

. poisson inju	uries XYZowne	d, exposure(n) irr			
Iteration 0: Iteration 1: Iteration 2:	Log likeliho Log likeliho Log likeliho	d = -23.027 d = -23.027 d = -23.027	197 177 177			
Poisson regre: Log likelihood	ssion 1 = -23.02717	7			Number of ob: LR chi2(1) Prob > chi2 Pseudo R2	s = 9 = 1.77 = 0.1836 = 0.0370
injuries	IRR	Std. err.	z	P> z	[95% conf.	interval]
XYZowned _cons ln(n)	1.463467 58.04416 1	.406872 8.558145 (exposure)	1.37 27.54	0.171 0.000	.8486578 43.47662	2.523675 77.49281

Note: _cons estimates baseline incidence rate.

We specified irr to see the IRRs rather than the underlying coefficients. We estimate that XYZ Airlines' injury rate is 1.46 times larger than that for other airlines, but the 95% confidence interval is 0.85 to 2.52; we cannot even reject the hypothesis that XYZ Airlines has a lower injury rate.

Technical note

In example 1, we assumed that each airline's exposure was proportional to its fraction of flights out of New York. What if "large" airlines, however, also used larger planes, and so had even more passengers than would be expected, given this measure of exposure? A better measure would be each airline's fraction of passengers on flights out of New York, a number that we do not have. Even so, we suppose that n represents this number to some extent, so a better estimate of the effect might be

. generate lnN=ln(n)						
. poisson inju	. poisson injuries XYZowned lnN					
<pre>Iteration 0: Log likelihood = -22.333875 Iteration 1: Log likelihood = -22.332276 Iteration 2: Log likelihood = -22.332276</pre>						
Poisson regree	ssion d = -22.332276				Number of ob: LR chi2(2) Prob > chi2 Pseudo R2	s = 9 = 19.15 = 0.0001 = 0.3001
injuries	Coefficient	Std. err.	z	P> z	[95% conf.	interval]
XYZowned lnN _cons	.6840667 1.424169 4.863891	.3895877 .3725155 .7090501	1.76 3.82 6.86	0.079 0.000 0.000	0795111 .6940517 3.474178	1.447645 2.154285 6.253603

Here rather than specifying the exposure() option, we explicitly included the variable that would normalize for exposure in the model. We did not specify the irr option, so we see coefficients rather than IRRs. We started with the model

rate =
$$e^{\beta_0 + \beta_1 XYZowned}$$

The observed counts are therefore

count =
$$ne^{\beta_0 + \beta_1 XYZowned} = e^{\ln(n) + \beta_0 + \beta_1 XYZowned}$$

which amounts to constraining the coefficient on ln(n) to 1. This is what was estimated when we specified the exposure (n) option. In the above model, we included the normalizing exposure ourselves and, rather than constraining the coefficient to be 1, estimated the coefficient.

The estimated coefficient is 1.42, a respectable distance away from 1, and is consistent with our speculation that larger airlines also use larger airplanes. With this small amount of data, however, we also have a wide confidence interval that includes 1.

Our estimated *coefficient* on XYZowned is now 0.684, and the implied IRR is $e^{0.684} \approx 1.98$ (which we could also see by typing poisson, irr). The 95% confidence interval for the coefficient still includes 0 (the interval for the IRR includes 1), so although the point estimate is now larger, we still cannot be certain of our results.

Our expert opinion would be that, although there is not enough evidence to support the charge, there is enough evidence to justify collecting more data.

Example 2

In a famous age-specific study of coronary disease deaths among male British doctors, Doll and Hill (1966) reported the following data (reprinted in Lash et al. [2021, 417]):

Smokers			Nonsmokers		
Age	Deaths	Person-years	Deaths	Person-years	
35-44	32	52,407	2	18,790	
45 - 54	104	43,248	12	10,673	
55 - 64	206	28,612	28	5,710	
65 - 74	186	12,663	28	2,585	
75 - 84	102	5,317	31	1,462	

The first step is to enter these data into Stata, which we have done:

```
. use https://www.stata-press.com/data/r19/dollhill3, clear (Doll and Hill (1966))
```

. list

	agecat	smokes	deaths	pyears
1. 2.	35-44 45-54	1	32 104	52,407 43,248
з.	55 - 64	1	206	28,612
4.	65-74	1	186	12,663
5.	75 - 84	1	102	5,317
6.	35-44	0	2	18,790
7.	45-54	0	12	10,673
8.	55-64	0	28	5,710
9.	65-74	0	28	2,585
0.	75 - 84	0	31	1,462

The most "natural" analysis of these data would begin by introducing indicator variables for each age category and one indicator for smoking:

. poisson deat	ths smokes i.	agecat, expo	sure(pyea	ars) irr		
Iteration 0: Iteration 1: Iteration 2: Iteration 3:	Log likeliho Log likeliho Log likeliho Log likeliho	od = -33.823 od = -33.600 od = -33.600 od = -33.600	284 9471 9153 9153			
Poisson regres	ssion 1 = -33.60015	3			Number of ob LR chi2(5) Prob > chi2 Pseudo R2	s = 10 = 922.93 = 0.0000 = 0.9321
deaths	IRR	Std. err.	z	P> z	[95% conf.	interval]
smokes	1.425519	.1530638	3.30	0.001	1.154984	1.759421
agecat 45-54 55-64 65-74 75-84	4.410584 13.8392 28.51678 40.45121	.8605197 2.542638 5.269878 7.775511	7.61 14.30 18.13 19.25	0.000 0.000 0.000 0.000	3.009011 9.654328 19.85177 27.75326	6.464997 19.83809 40.96395 58.95885
_cons ln(pyears)	.0003636 1	.0000697 (exposure)	-41.30	0.000	.0002497	.0005296

Note: _cons estimates baseline incidence rate.

In the above, we specified irr to obtain IRRs. We estimate that smokers have 1.43 times the mortality rate of nonsmokers. See, however, example 1 in [R] poisson postestimation.

Stored results

poisson stores the following in e():

```
Scalars
```

	e(N)	number of observations
	e(k)	number of parameters
	e(k_eq)	number of equations in e(b)
	e(k_eq_model)	number of equations in overall model test
	e(k_dv)	number of dependent variables
	e(df_m)	model degrees of freedom
	e(r2_p)	pseudo- R^2
	e(11)	log likelihood
	e(11_0)	log likelihood, constant-only model
	e(N_clust)	number of clusters
	e(chi2)	χ^2
	e(p)	<i>p</i> -value for model test
	e(rank)	rank of e(V)
	e(ic)	number of iterations
	e(rc)	return code
	e(converged)	1 if converged, 0 otherwise
Mac	ros	
	e(cmd)	poisson
	e(cmdline)	command as typed
	e(depvar)	name of dependent variable
	e(wtvpe)	weight type
		0 71

4

	e(wexp)	weight expression
	e(title)	title in estimation output
	e(clustvar)	name of cluster variable
	e(offset)	linear offset variable
	e(chi2type)	Wald or LR; type of model χ^2 test
	e(vce)	vcetype specified in vce()
	e(vcetype)	title used to label Std. err.
	e(opt)	type of optimization
	e(which)	max or min; whether optimizer is to perform maximization or minimization
	e(ml_method)	type of ml method
	e(user)	name of likelihood-evaluator program
	e(technique)	maximization technique
	e(properties)	b V
	e(estat_cmd)	program used to implement estat
	e(predict)	program used to implement predict
	e(asbalanced)	factor variables fvset as asbalanced
	e(asobserved)	factor variables fvset as asobserved
Mat	rices	
	e(b)	coefficient vector
	e(Cns)	constraints matrix
	e(ilog)	iteration log (up to 20 iterations)
	e(gradient)	gradient vector
	e(V)	variance-covariance matrix of the estimators
	e(V_modelbased)	model-based variance
Fund	ctions	
	e(sample)	marks estimation sample

In addition to the above, the following is stored in r():

Matrices

r(table) matrix containing the coefficients with their standard errors, test statistics, *p*-values, and confidence intervals

Note that results stored in r() are updated when the command is replayed and will be replaced when any r-class command is run after the estimation command.

Methods and formulas

The log likelihood (with weights w_j and offsets) is given by

$$\begin{split} \Pr(Y = y) &= \frac{e^{-\lambda}\lambda^y}{y!} \\ \xi_j &= \mathbf{x}_j \boldsymbol{\beta} + \mathrm{offset}_j \\ f(y_j) &= \frac{e^{-\exp(\xi_j)}e^{\xi_j y_j}}{y_j!} \\ \ln L &= \sum_{j=1}^n w_j \left\{ -e^{\xi_j} + \xi_j y_j - \ln(y_j!) \right\} \end{split}$$

This command supports the Huber/White/sandwich estimator of the variance and its clustered version using vce(robust) and vce(cluster *clustvar*), respectively. See [P] **_robust**, particularly *Maximum likelihood estimators* and *Methods and formulas*.

poisson also supports estimation with survey data. For details on VCEs with survey data, see [SVY] Variance estimation.

Siméon-Denis Poisson (1781–1840) was a French mathematician and physicist who contributed to several fields: his name is perpetuated in Poisson brackets, Poisson's constant, Poisson's differential equation, Poisson's integral, and Poisson's ratio. Among many other results, he produced a version of the law of large numbers. His rather misleadingly titled *Recherches sur la probabilité des jugements* embraces a complete treatise on probability, as the subtitle indicates, including what is now known as the Poisson distribution. That, however, was discovered earlier by the Huguenot–British mathematician Abraham de Moivre (1667–1754).

References

- Bru, B. 2001. "Siméon-Denis Poisson". In Statisticians of the Centuries, edited by C. C. Heyde and E. Seneta, 123–126. New York: Springer. https://doi.org/10.1007/978-1-4613-0179-0_25.
- Cameron, A. C., and P. K. Trivedi. 2013. Regression Analysis of Count Data. 2nd ed. New York: Cambridge University Press.

------. 2022. Microeconometrics Using Stata. 2nd ed. College Station, TX: Stata Press.

Chatterjee, S., and A. S. Hadi. 2012. Regression Analysis by Example. 5th ed. New York: Wiley.

- Clarke, R. D. 1946. An application of the Poisson distribution. *Journal of the Institute of Actuaries* 72: 481. https://doi. org/10.1017/S0020268100035435.
- Cochran, W. G. 1940. The analysis of variance when experimental errors follow the Poisson or binomial laws. Annals of Mathematical Statistics 11: 335–347. https://doi.org/10.1214/aoms/1177731871.

——. 1982. Contributions to Statistics. New York: Wiley.

Coleman, J. S. 1964. Introduction to Mathematical Sociology. New York: Free Press.

- Cummings, T. H., J. W. Hardin, A. C. McLain, J. R. Hussey, K. J. Bennett, and G. M. Wingood. 2015. Modeling heaped count data. Stata Journal 15: 457–479.
- Deb, P., E. C. Norton, and W. G. Manning. 2017. Health Econometrics Using Stata. College Station, TX: Stata Press.
- Doll, R., and A. B. Hill. 1966. Mortality of British doctors in relation to smoking: Observations on coronary thrombosis. Journal of the National Cancer Institute, Monographs 19: 205–268.
- Gould, W. W. 2011. Use poisson rather than regress; tell a friend. *The Stata Blog: Not Elsewhere Classified.* https://blog.stata.com/2011/08/22/use-poisson-rather-than-regress-tell-a-friend/.
- Harris, T., Z. Yang, and J. W. Hardin. 2012. Modeling underdispersed count data with generalized Poisson regression. *Stata Journal* 12: 736–747.
- Hilbe, J. M. 2014. Modeling Count Data. New York: Cambridge University Press.
- Jochmans, K., and V. Verardi. 2020. Fitting exponential regression models with two-way fixed effects. *Stata Journal* 20: 468–480.

Johnson, N. L., A. W. Kemp, and S. Kotz. 2005. Univariate Discrete Distributions. 3rd ed. New York: Wiley.

- Lash, T. L., T. J. VanderWeele, S. Haneuse, and K. J. Rothman. 2021. *Modern Epidemiology*. 4th ed. Philadelphia: Wolters Kluwer.
- Long, J. S. 1997. Regression Models for Categorical and Limited Dependent Variables. Thousand Oaks, CA: Sage.
- Long, J. S., and J. Freese. 2001. Predicted probabilities for count models. Stata Journal 1: 51-57.

McNeil, D. 1996. Epidemiological Research Methods. Chichester, UK: Wiley.

- Miranda, A., and S. Rabe-Hesketh. 2006. Maximum likelihood estimation of endogenous switching and sample selection models for binary, ordinal, and count variables. *Stata Journal* 6: 285–308.
- Newman, S. C. 2001. Biostatistical Methods in Epidemiology. New York: Wiley.
- Poisson, S. D. 1837. Recherches sur la probabilité des jugements en matière criminelle et en matière civile: précédées des règles générales du calcul des probabilités. Paris: Bachelier.
- Raciborski, R. 2011. Right-censored Poisson regression model. Stata Journal 11: 95-105.
- Rutherford, E., J. Chadwick, and C. D. Ellis. 1930. Radiations from Radioactive Substances. Cambridge: Cambridge University Press. https://doi.org/10.1017/CBO9780511707179.
- Rutherford, M. J., P. C. Lambert, and J. Thompson. 2010. Age-period-cohort modeling. Stata Journal 10: 606-627.
- Sasieni, P. D. 2012. Age-period-cohort models in Stata. Stata Journal 12: 45-60.
- Schonlau, M. 2005. Boosted regression (boosting): An introductory tutorial and a Stata plugin. Stata Journal 5: 330-354.
- Selvin, S. 2004. Statistical Analysis of Epidemiologic Data. 3rd ed. New York: Oxford University Press. https://doi.org/ 10.1093/acprof:oso/9780195172805.001.0001.
- Thorndike, F. 1926. Applications of Poisson's probability summation. *Bell System Technical Journal* 5: 604–624. https://doi.org/10.1002/j.1538-7305.1926.tb00126.x.
- von Bortkiewicz, L. 1898. Das Gesetz der Kleinen Zahlen. Leipzig: Teubner.
- Xu, X., and J. W. Hardin. 2016. Regression models for bivariate count outcomes. Stata Journal 16: 301-315.

Also see

- [R] **poisson postestimation** Postestimation tools for poisson
- [R] **glm** Generalized linear models
- [R] heckpoisson Poisson regression with sample selection
- [R] **nbreg** Negative binomial regression
- [R] npregress kernel Nonparametric kernel regression
- [R] npregress series Nonparametric series regression
- [R] tpoisson Truncated Poisson regression
- [R] **zip** Zero-inflated Poisson regression
- [BAYES] bayes: poisson Bayesian Poisson regression
- [FMM] fmm: poisson Finite mixtures of Poisson regression models
- [LASSO] Lasso intro Introduction to lasso
- [ME] mepoisson Multilevel mixed-effects Poisson regression
- [MI] Estimation Estimation commands for use with mi estimate
- [SVY] svy estimation Estimation commands for survey data
- [XT] xtpoisson Fixed-effects, random-effects, and population-averaged Poisson models

[U] 20 Estimation and postestimation commands

Stata, Stata Press, Mata, NetCourse, and NetCourseNow are registered trademarks of StataCorp LLC. Stata and Stata Press are registered trademarks with the World Intellectual Property Organization of the United Nations. StataNow is a trademark of StataCorp LLC. Other brand and product names are registered trademarks or trademarks of their respective companies. Copyright © 1985–2025 StataCorp LLC, College Station, TX, USA. All rights reserved.



For suggested citations, see the FAQ on citing Stata documentation.