

intreg — Interval regression[Description](#)
[Options](#)
[References](#)[Quick start](#)
[Remarks and examples](#)
[Also see](#)[Menu](#)
[Stored results](#)[Syntax](#)
[Methods and formulas](#)

Description

`intreg` fits a model in which the dependent variable may be measured as point data, interval data, left-censored data, or right-censored data. As such, it is a generalization of the models fit by `tobit`. The dependent variable must be specified using two *depvvars* that indicate how it was measured.

Quick start

Regression on `x1` and `x2` of an interval-measured dependent variable with lower endpoint `y_lower` and upper endpoint `y_upper`

```
intreg y_lower y_upper x1 x2
```

With robust standard errors

```
intreg y_lower y_upper x1 x2, vce(robust)
```

Model heteroskedasticity in the conditional variance as a function of `x3`

```
intreg y_lower y_upper x1 x2, het(x3)
```

Adjust for complex survey design using `svyset` data

```
svy: intreg y_lower y_upper x1 x2
```

Menu

Statistics > Linear models and related > Censored regression > Interval regression

Syntax

```
intreg depvar1 depvar2 [indepvars] [if] [in] [weight] [, options]
```

*depvar*₁ and *depvar*₂ should have the following form:

Type of data		<i>depvar</i> ₁	<i>depvar</i> ₂
point data	$a = [a, a]$	a	a
interval data	$[a, b]$	a	b
left-censored data	$(-\infty, b]$.	b
right-censored data	$[a, +\infty)$	a	.

<i>options</i>	Description
Model	
<code><u>noconstant</u></code>	suppress constant term
<code><u>het</u>(<i>varlist</i> [, <u>noconstant</u>])</code>	independent variables to model the variance; use <code>noconstant</code> to suppress constant term
<code><u>offset</u>(<i>varname</i>)</code>	include <i>varname</i> in model with coefficient constrained to 1
<code><u>constraints</u>(<i>constraints</i>)</code>	apply specified linear constraints
<code><u>collinear</u></code>	keep collinear variables
SE/Robust	
<code><u>vce</u>(<i>vcetype</i>)</code>	<i>vcetype</i> may be <code>oim</code> , <code>robust</code> , <code>cluster <i>clustvar</i></code> , <code>opg</code> , <code>bootstrap</code> , or <code>jackknife</code>
Reporting	
<code><u>level</u>(#)</code>	set confidence level; default is <code>level(95)</code>
<code><u>nocnsreport</u></code>	do not display constraints
<code><u>display_options</u></code>	control columns and column formats, row spacing, line width, display of omitted variables and base and empty cells, and factor-variable labeling
Maximization	
<code><u>maximize_options</u></code>	control the maximization process; seldom used
<code><u>coeflegend</u></code>	display legend instead of statistics

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

*depvar*₁, *depvar*₂, *indepvars*, and *varlist* may contain time-series operators; see [U] 11.4.4 Time-series varlists.

`bootstrap`, `by`, `fp`, `jackknife`, `mfp`, `nestreg`, `rolling`, `statsby`, `stepwise`, and `svy` are allowed; see [U] 11.1.10 Prefix commands.

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

`aweight`s are not allowed with the `jackknife` prefix; see [R] `jackknife`.

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

`aweight`s, `fweight`s, `iweight`s, and `pweight`s are allowed; see [U] 11.1.6 `weight`.

`coeflegend` does not appear in the dialog box.

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

Options

Model

`noconstant`; see [R] [estimation options](#).

`het(varlist [, noconstant])` specifies that *varlist* be included in the specification of the conditional variance. This *varlist* enters the variance specification collectively as multiplicative heteroskedasticity.

`offset(varname)`, `constraints(constraints)`, `collinear`; 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](#).

`nocnsreport`; see [R] [estimation options](#).

`display_options:` `nocl`, `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 option is available with `intreg` but is not shown in the dialog box:

`coeflegend`; see [R] [estimation options](#).

Remarks and examples

[stata.com](http://www.stata.com)

`intreg` is a generalization of the models fit by `tobit`. Cameron and Trivedi (2010, 548–550) discuss the differences among censored, truncated, and interval data. If you know that the value for the *j*th individual is somewhere in the interval $[y_{1j}, y_{2j}]$, then the likelihood contribution from this individual is simply $\Pr(y_{1j} \leq Y_j \leq y_{2j})$. For censored data, their likelihoods contain terms of the form $\Pr(Y_j \leq y_j)$ for left-censored data and $\Pr(Y_j \geq y_j)$ for right-censored data, where y_j is the observed censoring value and Y_j denotes the random variable representing the dependent variable in the model.

Hence, `intreg` can fit models for data where each observation represents interval data, left-censored data, right-censored data, or point data. Regardless of the type of observation, the data should be stored in the dataset as interval data; that is, two dependent variables, `depvar1` and `depvar2`, are used to hold the endpoints of the interval. If the data are left-censored, the lower endpoint is $-\infty$ and is represented by a missing value, `.`, or an extended missing value, `.a`, `.b`, `...`, `.z`, in `depvar1`.

If the data are right-censored, the upper endpoint is $+\infty$ and is represented by a missing value, . (or an extended missing value), in *depvar*₂. Point data are represented by the two endpoints being equal. Truly missing values of the dependent variable must be represented by missing values in both *depvar*₁ and *depvar*₂.

Interval data arise naturally in many contexts, such as wage data. Often you know only that, for example, a person's salary is between \$30,000 and \$40,000. Below we give an example for wage data and show how to set up *depvar*₁ and *depvar*₂.

► Example 1

We have a dataset that contains the yearly wages of working women. Women were asked via a questionnaire to indicate a category for their yearly income from employment. The categories were less than 5,000, 5,001–10,000, . . . , 25,001–30,000, 30,001–40,000, 40,001–50,000, and more than 50,000. The wage categories are stored in the *wagecat* variable.

```
. use http://www.stata-press.com/data/r14/womenwage
(Wages of women)
. tabulate wagecat
```

Wage category (\$1000s)	Freq.	Percent	Cum.
5	14	2.87	2.87
10	83	17.01	19.88
15	158	32.38	52.25
20	107	21.93	74.18
25	57	11.68	85.86
30	30	6.15	92.01
40	19	3.89	95.90
50	14	2.87	98.77
51	6	1.23	100.00
Total	488	100.00	

A value of 5 for *wagecat* represents the category less than 5,000, a value of 10 represents 5,001–10,000, . . . , and a value of 51 represents greater than 50,000.

To use *intreg*, we must create two variables, *wage1* and *wage2*, containing the lower and upper endpoints of the wage categories. Here is one way to do it. We first create a dataset containing the nine wage categories, lag the wage categories into *wage1*, and match-merge this dataset with nine observations back into the main one.

```
. by wagecat: keep if _n==1
(479 observations deleted)
. generate wage1 = wagecat[_n-1]
(1 missing value generated)
. keep wagecat wage1
. save lagwage
file lagwage.dta saved
. use http://www.stata-press.com/data/r14/womenwage
(Wages of women)
. merge m:1 wagecat using lagwage
```

Result	# of obs.
not matched	0
matched	488 (_merge==3)

Now we create the upper endpoint and list the new variables:

```
. generate wage2 = wagecat
. replace wage2 = . if wagecat == 51
(6 real changes made, 6 to missing)
. sort age, stable
. list wage1 wage2 in 1/10
```

	wage1	wage2
1.	.	5
2.	5	10
3.	5	10
4.	10	15
5.	.	5
6.	.	5
7.	.	5
8.	5	10
9.	5	10
10.	5	10

We can now run intreg:

```
. intreg wage1 wage2 age c.age#c.age nev_mar rural school tenure
Fitting constant-only model:
Iteration 0:   log likelihood = -967.24956
Iteration 1:   log likelihood = -967.1368
Iteration 2:   log likelihood = -967.1368
Fitting full model:
Iteration 0:   log likelihood = -856.65324
Iteration 1:   log likelihood = -856.33294
Iteration 2:   log likelihood = -856.33293
Interval regression           Number of obs   =           488
                             LR chi2(6)           =           221.61
                             Prob > chi2          =           0.0000
Log likelihood = -856.33293
```

	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
age	.7914438	.4433604	1.79	0.074	-.0775265	1.660414
c.age#c.age	-.0132624	.0073028	-1.82	0.069	-.0275757	.0010509
nev_mar	-.2075022	.8119581	-0.26	0.798	-1.798911	1.383906
rural	-3.043044	.7757324	-3.92	0.000	-4.563452	-1.522637
school	1.334721	.1357873	9.83	0.000	1.068583	1.600859
tenure	.8000664	.1045077	7.66	0.000	.5952351	1.004898
_cons	-12.70238	6.367117	-1.99	0.046	-25.1817	-.2230583
/lnsigma	1.987823	.0346543	57.36	0.000	1.919902	2.055744
sigma	7.299626	.2529634			6.82029	7.81265

```
14 left-censored observations
0 uncensored observations
6 right-censored observations
468 interval observations
```

We could also model these data by using an ordered probit model with `oprobit` (see [R] [oprobit](#)):

```
. oprobit wagecat age c.age#c.age nev_mar rural school tenure
Iteration 0:  log likelihood = -881.1491
Iteration 1:  log likelihood = -764.31729
Iteration 2:  log likelihood = -763.31191
Iteration 3:  log likelihood = -763.31049
Iteration 4:  log likelihood = -763.31049

Ordered probit regression                Number of obs   =       488
                                         LR chi2(6)      =       235.68
                                         Prob > chi2     =       0.0000
Log likelihood = -763.31049             Pseudo R2      =       0.1337
```

wagecat	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
age	.1674519	.0620333	2.70	0.007	.0458689	.289035
c.age#c.age	-.0027983	.0010214	-2.74	0.006	-.0048001	-.0007964
nev_mar	-.0046417	.1126737	-0.04	0.967	-.225478	.2161946
rural	-.5270036	.1100449	-4.79	0.000	-.7426875	-.3113196
school	.2010587	.0201189	9.99	0.000	.1616263	.2404911
tenure	.0989916	.0147887	6.69	0.000	.0700063	.127977
/cut1	2.650637	.8957245			.8950495	4.406225
/cut2	3.941018	.8979167			2.181134	5.700903
/cut3	5.085205	.9056582			3.310148	6.860263
/cut4	5.875534	.9120933			4.087864	7.663204
/cut5	6.468723	.918117			4.669247	8.268199
/cut6	6.922726	.9215455			5.11653	8.728922
/cut7	7.34471	.9237628			5.534168	9.155252
/cut8	7.963441	.9338881			6.133054	9.793828

We can directly compare the log likelihoods for the `intreg` and `oprobit` models because both likelihoods are discrete. If we had point data in our `intreg` estimation, the likelihood would be a mixture of discrete and continuous terms, and we could not compare it directly with the `oprobit` likelihood.

Here the `oprobit` log likelihood is significantly larger (that is, less negative), so it fits better than the `intreg` model. The `intreg` model assumes normality, but the distribution of wages is skewed and definitely nonnormal. Normality is more closely approximated if we model the log of wages.

```

. generate logwage1 = log(wage1)
(14 missing values generated)
. generate logwage2 = log(wage2)
(6 missing values generated)
. intreg logwage1 logwage2 age c.age#c.age nev_mar rural school tenure
Fitting constant-only model:
Iteration 0:   log likelihood = -889.23647
Iteration 1:   log likelihood = -889.06346
Iteration 2:   log likelihood = -889.06346
Fitting full model:
Iteration 0:   log likelihood = -773.81968
Iteration 1:   log likelihood = -773.36566
Iteration 2:   log likelihood = -773.36563
Interval regression               Number of obs   =       488
                                LR chi2(6)         =       231.40
                                Prob > chi2        =       0.0000
Log likelihood = -773.36563

```

	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
age	.0645589	.0249954	2.58	0.010	.0155689	.1135489
c.age#c.age	-.0010812	.0004115	-2.63	0.009	-.0018878	-.0002746
nev_mar	-.0058151	.0454867	-0.13	0.898	-.0949674	.0833371
rural	-.2098361	.0439454	-4.77	0.000	-.2959675	-.1237047
school	.0804832	.0076783	10.48	0.000	.0654341	.0955323
tenure	.0397144	.0058001	6.85	0.000	.0283464	.0510825
_cons	.7084023	.3593193	1.97	0.049	.0041495	1.412655
/lnsigma	-.906989	.0356265	-25.46	0.000	-.9768157	-.8371623
sigma	.4037381	.0143838			.3765081	.4329373

```

14 left-censored observations
0 uncensored observations
6 right-censored observations
468 interval observations

```

The log likelihood of this `intreg` model is close to the `oprobit` log likelihood, and the z statistics for both models are similar. ◀

□ Technical note

`intreg` has two parameterizations for the log-likelihood function: the transformed parameterization (β/σ , $1/\sigma$) and the untransformed parameterization (β , $\ln(\sigma)$). By default, the log likelihood for `intreg` is parameterized in the transformed parameter space. This parameterization tends to be more convergent, but it requires that any starting values and constraints have the same parameterization, and it prevents the estimation with multiplicative heteroskedasticity. Therefore, when the `het()` option is specified, `intreg` switches to the untransformed log likelihood for the fit of the conditional-variance model. Similarly, specifying `from()` or `constraints()` causes the optimization in the untransformed parameter space to allow constraints on (and starting values for) the coefficients on the covariates without reference to σ .

The estimation results are all stored in the $(\beta, \ln(\sigma))$ metric. □

Stored results

`intreg` stores the following in `e()`:

Scalars

<code>e(N)</code>	number of observations
<code>e(N_unc)</code>	number of uncensored observations
<code>e(N_ltc)</code>	number of left-censored observations
<code>e(N_rtc)</code>	number of right-censored observations
<code>e(N_int)</code>	number of interval observations
<code>e(k)</code>	number of parameters
<code>e(k_aux)</code>	number of auxiliary parameters
<code>e(k_eq)</code>	number of equations in <code>e(b)</code>
<code>e(k_eq_model)</code>	number of equations in overall model test
<code>e(k_dv)</code>	number of dependent variables
<code>e(df_m)</code>	model degrees of freedom
<code>e(ll)</code>	log likelihood
<code>e(ll_0)</code>	log likelihood, constant-only model
<code>e(ll_c)</code>	log likelihood, comparison model
<code>e(N_clust)</code>	number of clusters
<code>e(chi2)</code>	χ^2
<code>e(p)</code>	p -value for model χ^2 test
<code>e(sigma)</code>	sigma
<code>e(se_sigma)</code>	standard error of sigma
<code>e(rank)</code>	rank of <code>e(V)</code>
<code>e(rank0)</code>	rank of <code>e(V)</code> for constant-only model
<code>e(ic)</code>	number of iterations
<code>e(rc)</code>	return code
<code>e(converged)</code>	1 if converged, 0 otherwise

Macros

<code>e(cmd)</code>	<code>intreg</code>
<code>e(cmdline)</code>	command as typed
<code>e(depvar)</code>	names of dependent variables
<code>e(wtype)</code>	weight type
<code>e(wexp)</code>	weight expression
<code>e(title)</code>	title in estimation output
<code>e(clustvar)</code>	name of cluster variable
<code>e(offset)</code>	linear offset variable
<code>e(chi2type)</code>	Wald or LR; type of model χ^2 test
<code>e(vce)</code>	<code>vctype</code> specified in <code>vce()</code>
<code>e(vctype)</code>	title used to label Std. Err.
<code>e(het)</code>	heteroskedasticity, if <code>het()</code> specified
<code>e(ml_score)</code>	program used to implement <code>scores</code>
<code>e(opt)</code>	type of optimization
<code>e(which)</code>	max or min; whether optimizer is to perform maximization or minimization
<code>e(ml_method)</code>	type of ml method
<code>e(user)</code>	name of likelihood-evaluator program
<code>e(technique)</code>	maximization technique
<code>e(properties)</code>	<code>b V</code>
<code>e(predict)</code>	program used to implement <code>predict</code>
<code>e(footnote)</code>	program and arguments to display footnote
<code>e(marginsok)</code>	predictions allowed by <code>margins</code>
<code>e(asbalanced)</code>	factor variables <code>fvset</code> as <code>asbalanced</code>
<code>e(asobserved)</code>	factor variables <code>fvset</code> as <code>asobserved</code>

Matrices

<code>e(b)</code>	coefficient vector
<code>e(Cns)</code>	constraints matrix
<code>e(ilog)</code>	iteration log (up to 20 iterations)
<code>e(gradient)</code>	gradient vector
<code>e(V)</code>	variance-covariance matrix of the estimators
<code>e(V_modelbased)</code>	model-based variance

Functions

<code>e(sample)</code>	marks estimation sample
------------------------	-------------------------

Methods and formulas

See Wooldridge (2016, sec. 17.4) for an introduction to censored and truncated regression models.

The likelihood for `intreg` subsumes that of the `tobit` models.

Let $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ be the model. \mathbf{y} represents continuous outcomes—either observed or not observed. Our model assumes $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$.

For observations $j \in \mathcal{C}$, we observe y_j , that is, point data. Observations $j \in \mathcal{L}$ are left-censored; we know only that the unobserved y_j is less than or equal to $y_{\mathcal{L}j}$, a censoring value that we do know. Similarly, observations $j \in \mathcal{R}$ are right-censored; we know only that the unobserved y_j is greater than or equal to $y_{\mathcal{R}j}$. Observations $j \in \mathcal{I}$ are intervals; we know only that the unobserved y_j is in the interval $[y_{1j}, y_{2j}]$.

The log likelihood is

$$\begin{aligned} \ln L = & -\frac{1}{2} \sum_{j \in \mathcal{C}} w_j \left\{ \left(\frac{y_j - \mathbf{x}_j \boldsymbol{\beta}}{\sigma} \right)^2 + \log 2\pi\sigma^2 \right\} \\ & + \sum_{j \in \mathcal{L}} w_j \log \Phi \left(\frac{y_{\mathcal{L}j} - \mathbf{x}_j \boldsymbol{\beta}}{\sigma} \right) \\ & + \sum_{j \in \mathcal{R}} w_j \log \left\{ 1 - \Phi \left(\frac{y_{\mathcal{R}j} - \mathbf{x}_j \boldsymbol{\beta}}{\sigma} \right) \right\} \\ & + \sum_{j \in \mathcal{I}} w_j \log \left\{ \Phi \left(\frac{y_{2j} - \mathbf{x}_j \boldsymbol{\beta}}{\sigma} \right) - \Phi \left(\frac{y_{1j} - \mathbf{x}_j \boldsymbol{\beta}}{\sigma} \right) \right\} \end{aligned}$$

where $\Phi(\cdot)$ is the standard cumulative normal and w_j is the weight for the j th observation. If no weights are specified, $w_j = 1$. If `aweight`s are specified, $w_j = 1$, and σ is replaced by $\sigma/\sqrt{a_j}$ in the above, where a_j are the `aweight`s normalized to sum to N .

Maximization is as described in [R] `maximize`; the estimate reported as `_sigma` is $\hat{\sigma}$.

See Amemiya (1973) for a generalization of the tobit model to variable, but known, cutoffs.

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*.

`intreg` also supports estimation with survey data. For details on VCEs with survey data, see [SVY] `variance estimation`.

References

- Amemiya, T. 1973. Regression analysis when the dependent variable is truncated normal. *Econometrica* 41: 997–1016.
- Cameron, A. C., and P. K. Trivedi. 2010. *Microeconometrics Using Stata*. Rev. ed. College Station, TX: Stata Press.
- Conroy, R. M. 2005. *Stings in the tails: Detecting and dealing with censored data*. *Stata Journal* 5: 395–404.
- Davidson, R., and J. G. MacKinnon. 1993. *Estimation and Inference in Econometrics*. New York: Oxford University Press.
- Goldberger, A. S. 1983. Abnormal selection bias. In *Studies in Econometrics, Time Series, and Multivariate Statistics*, ed. S. Karlin, T. Amemiya, and L. A. Goodman, 67–84. New York: Academic Press.
- Hurd, M. 1979. Estimation in truncated samples when there is heteroscedasticity. *Journal of Econometrics* 11: 247–258.

- Long, J. S. 1997. *Regression Models for Categorical and Limited Dependent Variables*. Thousand Oaks, CA: Sage.
- Stewart, M. B. 1983. On least squares estimation when the dependent variable is grouped. *Review of Economic Studies* 50: 737–753.
- Wooldridge, J. M. 2016. *Introductory Econometrics: A Modern Approach*. 6th ed. Boston: Cengage.

Also see

- [R] **intreg postestimation** — Postestimation tools for intreg
- [R] **regress** — Linear regression
- [R] **tobit** — Tobit regression
- [SVY] **svy estimation** — Estimation commands for survey data
- [XT] **xtintreg** — Random-effects interval-data regression models
- [XT] **xttobit** — Random-effects tobit models
- [U] **20 Estimation and postestimation commands**