

**proportion** — Estimate proportions[Description](#)  
[Options](#)  
[References](#)[Quick start](#)  
[Remarks and examples](#)  
[Also see](#)[Menu](#)  
[Stored results](#)[Syntax](#)  
[Methods and formulas](#)

## Description

`proportion` produces estimates of proportions, along with standard errors, for the categories identified by the values in each variable of *varlist*.

## Quick start

Proportions, standard errors, and 95% CIs for each level of *v1*

```
proportion v1
```

Also compute statistics for *v2*

```
proportion v1 v2
```

Treat missing values of *v1* as a valid category

```
proportion v1, missing
```

As above, for each subpopulation defined by the levels of *catvar*

```
proportion v1, missing over(catvar)
```

Standardizing across strata defined by *svar* with stratum weight *wvar1*

```
proportion v1, stdize(svar) stdweight(wvar1)
```

Weighting by sampling weight *wvar2*

```
proportion v1 [pweight=wvar2]
```

## Menu

Statistics > Summaries, tables, and tests > Summary and descriptive statistics > Proportions

## Syntax

```
proportion varlist [if] [in] [weight] [, options]
```

<i>options</i>	Description
<b>Model</b>	
<code>stdize(<i>varname</i>)</code>	variable identifying strata for standardization
<code>stdweight(<i>varname</i>)</code>	weight variable for standardization
<code>nostdrescale</code>	do not rescale the standard weight variable
<code>no<del>l</del>abel</code>	suppress value labels from <i>varlist</i>
<code>missing</code>	treat missing values like other values
<b>if/in/over</b>	
<code>over(<i>varlist</i> [, <u>no<del>l</del>abel</u> ])</code>	group over subpopulations defined by <i>varlist</i> ; optionally, suppress group labels
<b>SE/Cluster</b>	
<code>vce(<i>vcetype</i>)</code>	<i>vcetype</i> may be <code>analytic</code> , <code>cluster <i>clustvar</i></code> , <code>bootstrap</code> , or <code>jackknife</code>
<b>Reporting</b>	
<code>level(#)</code>	set confidence level; default is <code>level(95)</code>
<code>citype(<i>citype</i>)</code>	method to compute limits of confidence intervals; default is <code>citype(logit)</code>
<code>noheader</code>	suppress table header
<code>nolegend</code>	suppress table legend
<code>display_options</code>	control column formats and line width
<code>coeflegend</code>	display legend instead of statistics

`bootstrap`, `jackknife`, `mi estimate`, `rolling`, `statsby`, and `svy` are allowed; see [U] 11.1.10 **Prefix commands**.

`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**.

`coeflegend` does not appear in the dialog box.

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

<i>citype</i>	Description
<code>logit</code>	calculate logit-transformed confidence intervals; the default
<code>agresti</code>	calculate Agresti–Coull confidence intervals
<code>exact</code>	calculate exact (Clopper–Pearson) confidence intervals
<code>jeffreys</code>	calculate Jeffreys confidence intervals
<code>normal</code>	calculate normal (Wald) confidence intervals
<code>wald</code>	synonym for <code>normal</code>
<code>wilson</code>	calculate Wilson confidence intervals

## Options

### Model

`stdize(varname)` specifies that the point estimates be adjusted by direct standardization across the strata identified by *varname*. This option requires the `stdweight()` option.

`stdweight(varname)` specifies the weight variable associated with the standard strata identified in the `stdize()` option. The standardization weights must be constant within the standard strata.

`nostdrescale` prevents the standardization weights from being rescaled within the `over()` groups. This option requires `stdize()` but is ignored if the `over()` option is not specified.

`nolabel` specifies that value labels attached to the variables in *varlist* be ignored.

`missing` specifies that missing values in *varlist* be treated as valid categories, rather than omitted from the analysis (the default).

### if/in/over

`over(varlist [ , nolabel ])` specifies that estimates be computed for multiple subpopulations, which are identified by the different values of the variables in *varlist*.

When this option is supplied with one variable name, such as `over(varname)`, the value labels of *varname* are used to identify the subpopulations. If *varname* does not have labeled values (or there are unlabeled values), the values themselves are used, provided that they are nonnegative integers. Noninteger values, negative values, and labels that are not valid Stata names are substituted with a default identifier.

When `over()` is supplied with multiple variable names, each subpopulation is assigned a unique default identifier.

`nolabel` requests that value labels attached to the variables identifying the subpopulations be ignored.

### SE/Cluster

`vce(vcetype)` specifies the type of standard error reported, which includes types that are derived from asymptotic theory (`analytic`), that allow for intragroup correlation (`cluster clustvar`), and that use bootstrap or jackknife methods (`bootstrap`, `jackknife`); see [R] [vce\\_option](#).

`vce(analytic)`, the default, uses the analytically derived variance estimator associated with the sample proportion.

### Reporting

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

`citype(citype)` specifies how to compute the limits of confidence intervals. *citype* may be one of `logit` (default), `agresti`, `exact`, `jeffreys`, `normal`, `wald`, or `wilson`.

`noheader` prevents the table header from being displayed. This option implies `nolegend`.

`nolegend` prevents the table legend identifying the subpopulations from being displayed.

*display\_options*: `cformat(%fmt)` and `nolstretch`; see [R] [estimation options](#).

The following option is available with `proportion` but is not shown in the dialog box:

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

## Remarks and examples

## ▷ Example 1

We can estimate the proportion of each repair rating in auto2.dta:

```
. use http://www.stata-press.com/data/r15/auto2
(1978 Automobile Data)
. proportion rep78
Proportion estimation          Number of obs   =          69
```

		Proportion	Std. Err.	Logit [95% Conf. Interval]	
rep78	Poor	.0289855	.0203446	.0070061	.1121326
	Fair	.115942	.0388245	.0580159	.2183014
	Average	.4347826	.0601159	.3207109	.556206
	Good	.2608696	.0532498	.1690271	.3798066
	Excellent	.1594203	.0443922	.089188	.2686455

Here we use the missing option to include missing values as a category of rep78:

```
. proportion rep78, missing
Proportion estimation          Number of obs   =          74
      _prop_6: rep78 = .
```

		Proportion	Std. Err.	Logit [95% Conf. Interval]	
rep78	Poor	.027027	.0189796	.0065484	.1047932
	Fair	.1081081	.0363433	.054094	.204402
	Average	.4054054	.0574637	.2977369	.523012
	Good	.2432432	.0502154	.1572724	.3563376
	Excellent	.1486486	.0416364	.0831005	.2517065
	_prop_6	.0675676	.0293776	.0278144	.1550743

▷ Example 2

We can also estimate proportions over groups:

```
. proportion rep78, over(foreign)
Proportion estimation      Number of obs   =           69
      Poor: rep78 = Poor
      Fair: rep78 = Fair
      Average: rep78 = Average
      Good: rep78 = Good
      Excellent: rep78 = Excellent
      Domestic: foreign = Domestic
      Foreign: foreign = Foreign
```

Over	Proportion	Std. Err.	Logit [95% Conf. Interval]	
Poor				
Domestic	.0416667	.0291477	.0100299	.1572433
Foreign	0	(no observations)		
Fair				
Domestic	.1666667	.0543607	.0839032	.3039797
Foreign	0	(no observations)		
Average				
Domestic	.5625	.0723605	.4169211	.6980553
Foreign	.1428571	.0782461	.0444941	.3736393
Good				
Domestic	.1875	.0569329	.0986718	.3272601
Foreign	.4285714	.1106567	.2333786	.6488451
Excellent				
Domestic	.0416667	.0291477	.0100299	.1572433
Foreign	.4285714	.1106567	.2333786	.6488451

## Stored results

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

### Scalars

<code>e(N)</code>	number of observations
<code>e(N_over)</code>	number of subpopulations
<code>e(N_stdize)</code>	number of standard strata
<code>e(N_clust)</code>	number of clusters
<code>e(k_eq)</code>	number of equations in <code>e(b)</code>
<code>e(df_r)</code>	sample degrees of freedom
<code>e(rank)</code>	rank of <code>e(V)</code>

### Macros

<code>e(cmd)</code>	<code>proportion</code>
<code>e(cmdline)</code>	command as typed
<code>e(varlist)</code>	<i>varlist</i>
<code>e(stdize)</code>	<i>varname</i> from <code>stdize()</code>
<code>e(stdweight)</code>	<i>varname</i> from <code>stdweight()</code>
<code>e(wtype)</code>	weight type
<code>e(wexp)</code>	weight expression
<code>e(title)</code>	title in estimation output
<code>e(cluster)</code>	name of cluster variable
<code>e(over)</code>	<i>varlist</i> from <code>over()</code>
<code>e(over_labels)</code>	labels from <code>over()</code> variables
<code>e(over_namelist)</code>	names from <code>e(over_labels)</code>
<code>e(namelist)</code>	<code>proportion</code> identifiers
<code>e(label#)</code>	labels from <i>#</i> th variable in <i>varlist</i>
<code>e(vce)</code>	<i>vcetype</i> specified in <code>vce()</code>
<code>e(vcetype)</code>	title used to label Std. Err.
<code>e(properties)</code>	<code>b V</code>
<code>e(estat_cmd)</code>	program used to implement <code>estat</code>
<code>e(marginsnotok)</code>	predictions disallowed by <code>margins</code>

### Matrices

<code>e(b)</code>	vector of proportion estimates
<code>e(V)</code>	(co)variance estimates
<code>e(_N)</code>	vector of numbers of nonmissing observations
<code>e(_N_stdsum)</code>	number of nonmissing observations within the standard strata
<code>e(_p_stdize)</code>	standardizing proportions
<code>e(error)</code>	error code corresponding to <code>e(b)</code>

### Functions

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

## Methods and formulas

Proportions are means of indicator variables; see [R] [mean](#).

## Confidence intervals

For an overview of confidence interval methods for binomial proportions, see [Dean and Pagano \(2015\)](#).

Given  $k$  successes of  $n$  trials, the estimated proportion (probability of a success) is  $\hat{p} = k/n$  with estimated standard error  $\hat{s} = \sqrt{\hat{p}(1-\hat{p})/n}$ .

The logit-transformed confidence interval is given by

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) \pm t_{1-\alpha/2,\nu} \frac{\hat{s}}{\hat{p}(1-\hat{p})}$$

where  $t_{p,\nu}$  is the  $p$ th quantile of Student's  $t$  distribution with  $\nu$  degrees of freedom.

The endpoints of this confidence interval are transformed back to the proportion metric by using the inverse of the logit transform

$$f^{-1}(y) = \frac{e^y}{1 + e^y}$$

Hence, the displayed confidence intervals for proportions are

$$f^{-1} \left\{ \ln \left( \frac{\hat{p}}{1 - \hat{p}} \right) \pm t_{1-\alpha/2, \nu} \frac{\hat{s}}{\hat{p}(1 - \hat{p})} \right\}$$

The Wald-type  $100(1 - \alpha)\%$  confidence interval is given by

$$\hat{p} \pm t_{1-\alpha/2, \nu} \hat{s}$$

The Wilson interval is given by

$$\frac{\hat{p} + z_{1-\alpha/2}^2/2n \pm z_{1-\alpha/2} \sqrt{\hat{s} + z_{1-\alpha/2}^2/4n^2}}{1 + z_{1-\alpha/2}^2/n}$$

where  $z_p$  is the  $p$ th quantile of the standard normal distribution.

The exact (Clopper–Pearson) interval is given by

$$\left\{ \hat{p} - \frac{\nu_1 F_{\alpha/2, \nu_1, \nu_2}}{\nu_2 + \nu_1 F_{\alpha/2, \nu_1, \nu_2}}; \hat{p} + \frac{\nu_3 F_{\alpha/2, \nu_3, \nu_4}}{\nu_4 + \nu_3 F_{\alpha/2, \nu_3, \nu_4}} \right\}$$

where  $\nu_1 = 2k$ ,  $\nu_2 = 2(n - k + 1)$ ,  $\nu_3 = 2(k + 1)$ ,  $\nu_4 = 2(n - k)$ , and  $F_{p, \nu_1, \nu_2}$  is the  $p$ th quantile of an  $F$  distribution with  $\nu_1$  and  $\nu_2$  degrees of freedom.

The Jeffreys interval is given by

$$\{\hat{p} - \text{Beta}_{\alpha/2, \alpha_1, \beta_2}; \hat{p} + \text{Beta}_{1-\alpha/2, \alpha_1, \beta_2}\}$$

where  $\alpha_1 = k + 0.5$ ,  $\beta_1 = n - k + 0.5$ , and  $\text{Beta}_{p, \alpha_1, \beta_2}$  is the  $p$ th quantile of a Beta distribution with  $\alpha_1$  and  $\beta_1$  degrees of freedom.

The Agresti–Coull interval is given by

$$\tilde{p} \pm z_{1-\alpha/2} \sqrt{\tilde{p}(1 - \tilde{p})/\tilde{n}}$$

where  $\tilde{k} = k + z_{1-\alpha/2}^2$ ,  $\tilde{n} = n + 2z_{1-\alpha/2}^2$ , and  $\tilde{p} = \tilde{k}/\tilde{n}$ .

When degrees of freedom  $\nu$  are posted to `e(df_r)`, the Wilson, exact, Jeffreys, and Agresti–Coull intervals use  $n^*$  in place of  $n$ , where

$$n^* = \frac{\hat{p}(1 - \hat{p})}{\hat{s}^2} \left\{ \frac{z_{1-\alpha/2}}{t_{1-\alpha/2, \nu}} \right\}^2$$

## References

- Cochran, W. G. 1977. *Sampling Techniques*. 3rd ed. New York: Wiley.
- Dean, N., and M. Pagano. 2015. Evaluating confidence interval methods for binomial proportions in clustered surveys. *Journal of Survey Statistics and Methodology* 3: 484–503.
- Stuart, A., and J. K. Ord. 1994. *Kendall's Advanced Theory of Statistics: Distribution Theory, Vol I*. 6th ed. London: Arnold.

## Also see

- [R] **proportion postestimation** — Postestimation tools for proportion
- [R] **mean** — Estimate means
- [R] **ratio** — Estimate ratios
- [R] **total** — Estimate totals
- [MI] **estimation** — Estimation commands for use with mi estimate
- [SVY] **direct standardization** — Direct standardization of means, proportions, and ratios
- [SVY] **poststratification** — Poststratification for survey data
- [SVY] **subpopulation estimation** — Subpopulation estimation for survey data
- [SVY] **svy estimation** — Estimation commands for survey data
- [SVY] **variance estimation** — Variance estimation for survey data
- [U] **20 Estimation and postestimation commands**