

**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
```

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

```
proportion v1 v2, 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
<b>if/in/over</b>	
<code>over(<i>varlist<sub>o</sub></i>)</code>	group over subpopulations defined by <i>varlist<sub>o</sub></i>
<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>percent</code>	report estimated proportions as percentages
<code>noheader</code>	suppress table header
<code>display_options</code>	control column formats, row spacing, line width, display of omitted variables and base and empty cells, and factor-variable labeling
<code>coeflegend</code>	display legend instead of statistics

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

Only numeric, nonnegative, integer-valued variables are allowed in *varlist*.

`bootstrap`, `collect`, `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.

### if/in/over

`over(varlisto)` specifies that estimates be computed for multiple subpopulations, which are identified by the different values of the variables in *varlist<sub>o</sub>*. Only numeric, nonnegative, integer-valued variables are allowed in `over(varlisto)`.

### 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`.

`percent` specifies that the proportions be reported as percentages.

`noheader` prevents the table header from being displayed.

*display\_options*: `noomitted`, `vsquish`, `noemptycells`, `baselevels`, `allbaselevels`, `nofvlabel`, `fvwrap(#)`, `fvwrapon(style)`, `cformat(%fmt)`, and `no1stretch`; 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 https://www.stata-press.com/data/r17/auto2
(1978 automobile data)
```

```
. proportion rep78
```

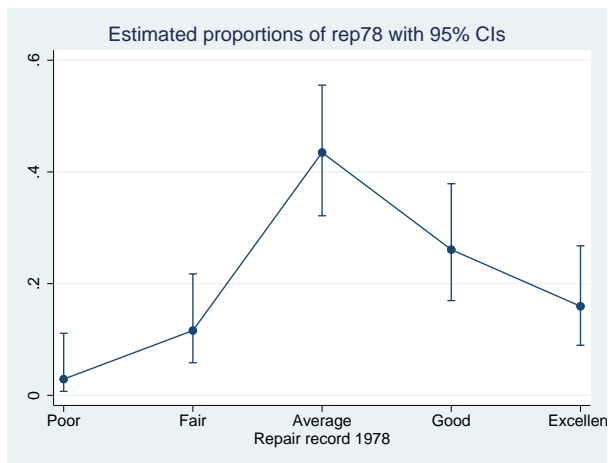
Proportion estimation Number of obs = 69

	Proportion	Std. err.	Logit	
			[95% conf. interval]	
rep78				
Poor	.0289855	.0201966	.0070794	.1110924
Fair	.115942	.0385422	.058317	.2173648
Average	.4347826	.0596787	.3214848	.5553295
Good	.2608696	.0528625	.1695907	.3788629
Excellent	.1594203	.0440694	.0895793	.267702

`marginsplot` will produce a graph of the results from `proportion`:

```
. marginsplot
```

Variables that uniquely identify proportions: **rep78**



► Example 2

We can also estimate proportions over groups:

```
. proportion rep78, over(foreign)
Proportion estimation                               Number of obs = 69
```

	Proportion	Std. err.	Logit [95% conf. interval]	
rep78@foreign				
Poor Domestic	.0416667	.0288424	.0101825	.1552326
Poor Foreign	0	(no observations)		
Fair Domestic	.1666667	.0537914	.084534	.3022522
Fair Foreign	0	(no observations)		
Average Domestic	.5625	.0716027	.4184154	.6967587
Average Foreign	.1428571	.0763604	.0458191	.3664757
Good Domestic	.1875	.0563367	.0993684	.3255432
Good Foreign	.4285714	.1079898	.2372889	.6438783
Excellent Domestic	.0416667	.0288424	.0101825	.1552326
Excellent Foreign	.4285714	.1079898	.2372889	.6438783

To see the results as percentages instead of proportions, we add the percent option:

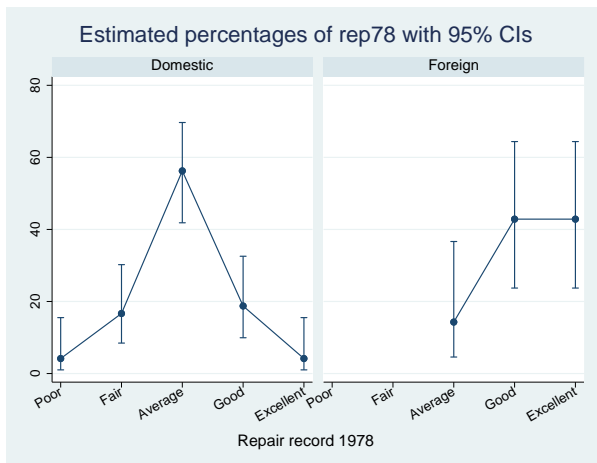
```
. proportion rep78, over(foreign) percent
Percent estimation                               Number of obs = 69
```

	Percent	Std. err.	Logit [95% conf. interval]	
rep78@foreign				
Poor Domestic	4.17	2.88	1.02	15.52
Poor Foreign	0.00	(no observations)		
Fair Domestic	16.67	5.38	8.45	30.23
Fair Foreign	0.00	(no observations)		
Average Domestic	56.25	7.16	41.84	69.68
Average Foreign	14.29	7.64	4.58	36.65
Good Domestic	18.75	5.63	9.94	32.55
Good Foreign	42.86	10.80	23.73	64.39
Excellent Domestic	4.17	2.88	1.02	15.52
Excellent Foreign	42.86	10.80	23.73	64.39

We can now use marginsplot to graph the percentages for each group. We add the bydimension(foreign) option to plot the groups in separate graph panels. The xlabel(, angle(30)) option prevents the x-axis labels from running into each other.

## 6 proportion — Estimate proportions

```
. marginsplot, bydimension(foreign) xlabel(, angle(30))
Variables that uniquely identify proportions: rep78 foreign
```



We estimate that only 19% of domestic cars have good repair records and only 4% have excellent repair records. For foreign cars, however, we find that 43% have good repair records and 43% have excellent repair records.



### Example 3

Instead of estimating percentages within the foreign and domestic groupings, we might want to know overall percentages. For instance, what percentage of all cars are foreign and have excellent repair records? What percentage are domestic and have average records? We can obtain all such percentages by specifying an interaction between `rep78` and `foreign`.

```
. proportion rep78#foreign, percent
```

Percent estimation

Number of obs = 69

	Percent	Std. err.	Logit [95% conf. interval]	
rep78#foreign				
Poor#Domestic	2.90	2.02	0.71	11.11
Poor#Foreign	0.00	(no observations)		
Fair#Domestic	11.59	3.85	5.83	21.74
Fair#Foreign	0.00	(no observations)		
Average#Domestic	39.13	5.88	28.21	51.26
Average#Foreign	4.35	2.46	1.38	12.86
Good#Domestic	13.04	4.05	6.85	23.44
Good#Foreign	13.04	4.05	6.85	23.44
Excellent#Domestic	2.90	2.02	0.71	11.11
Excellent#Foreign	13.04	4.05	6.85	23.44

Looking at the last line of this output, we estimate that 13% of all cars are foreign with excellent repair records.



## 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(clustvar)</code>	name of cluster variable
<code>e(over)</code>	<i>varlist</i> from <code>over()</code>
<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(freq)</code>	vector of frequency estimates
<code>e(error)</code>	error code corresponding to <code>e(b)</code>

### Functions

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

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

### Matrices

<code>r(table)</code>	matrix containing the coefficients with their standard errors, test statistics, <i>p</i> -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

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  $\tilde{k} = k + z_{1-\alpha/2}^2/2$ ,  $\tilde{n} = n + z_{1-\alpha/2}^2$ , 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. <https://doi.org/10.1093/jssam/smv024>.
- 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**