

Maximization

<code>irls</code>	use iterated, reweighted least-squares optimization; the default
<code>ml</code>	use maximum likelihood optimization
<code>maximize_options</code>	control the maximization process; seldom used
<code>fisher(#)</code>	Fisher scoring steps
<code>search</code>	search for good starting values
<code>coeflegend</code>	display legend instead of statistics

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

`deprvar` and `indepvars` may contain time-series operators; see [U] 11.4.4 **Time-series varlists**.

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

`vce(bootstrap)`, `vce(jackknife)`, and `vce(jackknife1)` are not allowed with the `mi estimate` prefix; see [MI] **mi estimate**.

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

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

`fweight`s, `aweight`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.

Menu

Statistics > Generalized linear models > GLM for the binomial family

Description

`binreg` fits generalized linear models for the binomial family. It estimates odds ratios, risk ratios, health ratios, and risk differences. The available links are

Option	Implied link	Parameter
<code>or</code>	<code>logit</code>	odds ratios = $\exp(\beta)$
<code>rr</code>	<code>log</code>	risk ratios = $\exp(\beta)$
<code>hr</code>	<code>log complement</code>	health ratios = $\exp(\beta)$
<code>rd</code>	<code>identity</code>	risk differences = β

Estimates of odds, risk, and health ratios are obtained by exponentiating the appropriate coefficients. The `or` option produces the same results as Stata's `logistic` command, and `or coefficients` yields the same results as the `logit` command. When no link is specified, `or` is assumed.

Options

Model

`noconstant`; see [R] **estimation options**.

`or` requests the logit link and results in odds ratios if `coefficients` is not specified.

`rr` requests the log link and results in risk ratios if `coefficients` is not specified.

`hr` requests the log-complement link and results in health ratios if `coefficients` is not specified.

`rd` requests the identity link and results in risk differences.

`n(#|varname)` specifies either a constant integer to use as the denominator for the binomial family or a variable that holds the denominator for each observation.

`exposure(varname)`, `offset(varname)`, `constraints(constraints)`, `collinear`; see [R] [estimation options](#). `constraints(constraints)` and `collinear` are not allowed with `irls`.

`mu(varname)` specifies `varname` containing an initial estimate for the mean of `depvar`. This option can be useful if you encounter convergence difficulties. `init(varname)` is a synonym.

SE/Robust

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

`vce(eim)`, the default, uses the expected information matrix (EIM) for the variance estimator.

`binreg` also allows the following:

`vce(hac kernel [#])` specifies that a heteroskedasticity- and autocorrelation-consistent (HAC) variance estimate be used. HAC refers to the general form for combining weighted matrices to form the variance estimate. There are three kernels built into `binreg`. `kernel` is a user-written program or one of

`nwest` | `gallant` | `anderson`

If `#` is not specified, $N - 2$ is assumed.

`vce(jackknife1)` specifies that the one-step jackknife estimate of variance be used.

`vce(unbiased)` specifies that the unbiased sandwich estimate of variance be used.

`t(varname)` specifies the variable name corresponding to time; see [TS] [tsset](#). `binreg` does not always need to know `t()`, though it does if `vce(hac ...)` is specified. Then you can either specify the time variable with `t()`, or you can `tsset` your data before calling `binreg`. When the time variable is required, `binreg` assumes that the observations are spaced equally over time.

`vfactor(#)` specifies a scalar by which to multiply the resulting variance matrix. This option allows users to match output with other packages, which may apply degrees of freedom or other small-sample corrections to estimates of variance.

`disp(#)` multiplies the variance of `depvar` by `#` and divides the deviance by `#`. The resulting distributions are members of the quasilielihood family.

`scale(x2|dev|#)` overrides the default scale parameter. This option is allowed only with Hessian (information matrix) variance estimates.

By default, `scale(1)` is assumed for the discrete distributions (binomial, Poisson, and negative binomial), and `scale(x2)` is assumed for the continuous distributions (Gaussian, gamma, and inverse Gaussian).

`scale(x2)` specifies that the scale parameter be set to the Pearson chi-squared (or generalized chi-squared) statistic divided by the residual degrees of freedom, which was recommended by [McCullagh and Nelder \(1989\)](#) as a good general choice for continuous distributions.

`scale(dev)` sets the scale parameter to the deviance divided by the residual degrees of freedom. This option provides an alternative to `scale(x2)` for continuous distributions and overdispersed or underdispersed discrete distributions.

`scale(#)` sets the scale parameter to `#`.

Reporting

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

`coefficients` displays the nonexponentiated coefficients and corresponding standard errors and confidence intervals. This option has no effect when the `rd` option is specified, because it always presents the nonexponentiated coefficients.

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

`display_options`: `noomitted`, `vsquish`, `noemptycells`, `baselevels`, `allbaselevels`, `nofvlabel`, `fvwrap(#)`, `fvwrapon(style)`, `cformat(%fmt)`, `pformat(%fmt)`, `sformat(%fmt)`, and `no1stretch`; see [R] [estimation options](#).

Maximization

`irls` requests iterated, reweighted least-squares (IRLS) optimization of the deviance instead of Newton–Raphson optimization of the log likelihood. This option is the default.

`ml` requests that optimization be carried out by using Stata’s `ml` command; see [R] [ml](#).

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

Setting the optimization method to `ml`, with `technique()` set to something other than `BHHH`, changes the `vcetype` to `vce(oim)`. Specifying `technique(bhhh)` changes `vcetype` to `vce(opg)`.

`fisher(#)` specifies the number of Newton–Raphson steps that should use the Fisher scoring Hessian or EIM before switching to the observed information matrix (OIM). This option is available only if `ml` is specified and is useful only for Newton–Raphson optimization.

`search` specifies that the command search for good starting values. This option is available only if `ml` is specified and is useful only for Newton–Raphson optimization.

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

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

Remarks and examples[stata.com](http://www.stata.com)

Wacholder (1986) suggests methods for estimating risk ratios and risk differences from prospective binomial data. These estimates are obtained by selecting the proper link functions in the generalized linear-model framework. (See [Methods and formulas](#) for details; also see [R] [glm](#).)

▷ Example 1

Wacholder (1986) presents an example, using data from Wright et al. (1983), of an investigation of the relationship between alcohol consumption and the risk of a low-birthweight baby. Covariates examined included whether the mother smoked (yes or no), mother’s social class (three levels), and drinking frequency (light, moderate, or heavy). The data for the 18 possible categories determined by the covariates are illustrated below.

Let’s first describe the data and list a few observations.

```
. use http://www.stata-press.com/data/r13/binreg
. list
```

	category	n_lbw_~s	n_women	alcohol	smokes	social
1.	1	11	84	heavy	nonsmoker	1
2.	2	5	79	moderate	nonsmoker	1
3.	3	11	169	light	nonsmoker	1
4.	4	6	28	heavy	smoker	1
5.	5	3	13	moderate	smoker	1
6.	6	1	26	light	smoker	1
7.	7	4	22	heavy	nonsmoker	2
8.	8	3	25	moderate	nonsmoker	2
9.	9	12	162	light	nonsmoker	2
10.	10	4	17	heavy	smoker	2
11.	11	2	7	moderate	smoker	2
12.	12	6	38	light	smoker	2
13.	13	0	14	heavy	nonsmoker	3
14.	14	1	18	moderate	nonsmoker	3
15.	15	12	91	light	nonsmoker	3
16.	16	7	19	heavy	smoker	3
17.	17	2	18	moderate	smoker	3
18.	18	8	70	light	smoker	3

Each observation corresponds to one of the 18 covariate structures. The number of low-birthweight babies from n_women in each category is given by the n_lbw_babies variable.

We begin by estimating risk ratios:

```
. binreg n_lbw_babies i.soc i.alc i.smo, n(n_women) rr
Iteration 1: deviance = 14.2879
Iteration 2: deviance = 13.607
Iteration 3: deviance = 13.60503
Iteration 4: deviance = 13.60503

Generalized linear models                               No. of obs   =       18
Optimization      : MQL Fisher scoring                 Residual df   =       12
                   (IRLS EIM)                         Scale parameter =       1
Deviance          = 13.6050268                        (1/df) Deviance = 1.133752
Pearson          = 11.51517095                        (1/df) Pearson  = .9595976

Variance function: V(u) = u*(1-u/n_women)             [Binomial]
Link function     : g(u) = ln(u/n_women)              [Log]
BIC                                                        = -21.07943
```

n_lbw_babies	EIM					
	Risk Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
social						
2	1.340001	.3127382	1.25	0.210	.848098	2.11721
3	1.349487	.3291488	1.23	0.219	.8366715	2.176619
alcohol						
moderate	1.191157	.3265354	0.64	0.523	.6960276	2.038503
heavy	1.974078	.4261751	3.15	0.002	1.293011	3.013884
smokes						
smoker	1.648444	.332875	2.48	0.013	1.109657	2.448836
_cons	.0630341	.0128061	-13.61	0.000	.0423297	.0938656

By default, Stata reports the risk ratios (the exponentiated regression coefficients) estimated by the model. We can see that the risk ratio comparing heavy drinkers with light drinkers, after adjusting for smoking and social class, is 1.974078. That is, mothers who drink heavily during their pregnancy have approximately twice the risk of delivering low-birthweight babies as mothers who are light drinkers.

The nonexponentiated coefficients can be obtained with the `coefficients` option:

```
. binreg n_lbw_babies i.soc i.alc i.smo, n(n_women) rr coefficients
Iteration 1:  deviance =   14.2879
Iteration 2:  deviance =   13.607
Iteration 3:  deviance =  13.60503
Iteration 4:  deviance =  13.60503

Generalized linear models                               No. of obs   =       18
Optimization   : MQL Fisher scoring                    Residual df   =       12
                 (IRLS EIM)                           Scale parameter =       1
Deviance       =  13.6050268                          (1/df) Deviance =  1.133752
Pearson        =  11.51517095                          (1/df) Pearson  =  .9595976

Variance function: V(u) = u*(1-u/n_women)             [Binomial]
Link function    : g(u) = ln(u/n_women)                [Log]
BIC                                                       = -21.07943
```

n_lbw_babies	EIM		z	P> z	[95% Conf. Interval]	
	Coef.	Std. Err.				
social						
2	.2926702	.2333866	1.25	0.210	-.1647591	.7500994
3	.2997244	.2439066	1.23	0.219	-.1783238	.7777726
alcohol						
moderate	.1749248	.274133	0.64	0.523	-.362366	.7122156
heavy	.6801017	.2158856	3.15	0.002	.2569737	1.10323
smokes						
smoker	.4998317	.2019329	2.48	0.013	.1040505	.8956129
_cons	-2.764079	.2031606	-13.61	0.000	-3.162266	-2.365891

Risk differences are obtained with the rd option:

```
. binreg n_lbw_babies i.soc i.alc i.smo, n(n_women) rd
Iteration 1: deviance = 18.67277
Iteration 2: deviance = 14.94364
Iteration 3: deviance = 14.9185
Iteration 4: deviance = 14.91762
Iteration 5: deviance = 14.91758
Iteration 6: deviance = 14.91758
Iteration 7: deviance = 14.91758

Generalized linear models                               No. of obs   =       18
Optimization      : MQL Fisher scoring                 Residual df   =       12
                   (IRLS EIM)                       Scale parameter =       1
Deviance          = 14.91758277                       (1/df) Deviance = 1.243132
Pearson          = 12.60353235                       (1/df) Pearson = 1.050294

Variance function: V(u) = u*(1-u/n_women)             [Binomial]
Link function     : g(u) = u/n_women                  [Identity]
BIC               = -19.76688
```

n_lbw_babies	EIM		z	P> z	[95% Conf. Interval]	
	Risk Diff.	Std. Err.				
social						
2	.0263817	.0232124	1.14	0.256	-.0191137	.0718771
3	.0365553	.0268668	1.36	0.174	-.0161026	.0892132
alcohol						
moderate	.0122539	.0257713	0.48	0.634	-.0382569	.0627647
heavy	.0801291	.0302878	2.65	0.008	.020766	.1394921
smokes						
smoker	.0542415	.0270838	2.00	0.045	.0011582	.1073248
_cons	.059028	.0160693	3.67	0.000	.0275327	.0905232

The risk difference between heavy drinkers and light drinkers is 0.0801291. Because the risk differences are obtained directly from the coefficients estimated by using the identity link, the `coefficients` option has no effect here.

Health ratios are obtained with the `hr` option. The health ratios (exponentiated coefficients for the log-complement link) are reported directly.

```
. binreg n_lbwbabies i.soc i.alc i.smo, n(n_women) hr
Iteration 1: deviance = 21.15233
Iteration 2: deviance = 15.16467
Iteration 3: deviance = 15.13205
Iteration 4: deviance = 15.13114
Iteration 5: deviance = 15.13111
Iteration 6: deviance = 15.13111
Iteration 7: deviance = 15.13111

Generalized linear models                No. of obs    =      18
Optimization   : MQL Fisher scoring      Residual df   =      12
                  (IRLS EIM)            Scale parameter =      1
Deviance       = 15.13110545             (1/df) Deviance = 1.260925
Pearson        = 12.84203917             (1/df) Pearson  = 1.07017

Variance function: V(u) = u*(1-u/n_women) [Binomial]
Link function    : g(u) = ln(1-u/n_women) [Log complement]
BIC              = -19.55336
```

n_lbwbabies	EIM			z	P> z	[95% Conf. Interval]	
	HR	Std. Err.					
social							
2	.9720541	.024858	-1.11	0.268	.9245342	1.022017	
3	.9597182	.0290412	-1.36	0.174	.9044535	1.01836	
alcohol							
moderate	.9871517	.0278852	-0.46	0.647	.9339831	1.043347	
heavy	.9134243	.0325726	-2.54	0.011	.8517631	.9795493	
smokes							
smoker	.9409983	.0296125	-1.93	0.053	.8847125	1.000865	
_cons	.9409945	.0163084	-3.51	0.000	.9095674	.9735075	

(HR) Health ratios

To see the nonexponentiated coefficients, we can specify the `coefficients` option.

Stored results

`binreg`, `irls` stores the following in `e()`:

Scalars

<code>e(N)</code>	number of observations
<code>e(k)</code>	number of parameters
<code>e(k_eq_model)</code>	number of equations in overall model test
<code>e(df_m)</code>	model degrees of freedom
<code>e(df)</code>	residual degrees of freedom
<code>e(phi)</code>	model scale parameter
<code>e(disp)</code>	dispersion parameter
<code>e(bic)</code>	model BIC
<code>e(N_clust)</code>	number of clusters
<code>e(deviance)</code>	deviance
<code>e(deviance_s)</code>	scaled deviance
<code>e(deviance_p)</code>	Pearson deviance
<code>e(deviance_ps)</code>	scaled Pearson deviance
<code>e(dispers)</code>	dispersion
<code>e(dispers_s)</code>	scaled dispersion
<code>e(dispers_p)</code>	Pearson dispersion
<code>e(dispers_ps)</code>	scaled Pearson dispersion
<code>e(vf)</code>	factor set by <code>vfactor()</code> , 1 if not set
<code>e(rank)</code>	rank of $e(V)$
<code>e(rc)</code>	return code

Macros

<code>e(cmd)</code>	<code>binreg</code>
<code>e(cmdline)</code>	command as typed
<code>e(depvar)</code>	name of dependent variable
<code>e(eform)</code>	<code>eform()</code> option implied by <code>or</code> , <code>rr</code> , <code>hr</code> , or <code>rd</code>
<code>e(varfunc)</code>	program to calculate variance function
<code>e(varfunct)</code>	variance title
<code>e(varfunctf)</code>	variance function
<code>e(link)</code>	program to calculate link function
<code>e(linkt)</code>	link title
<code>e(linkf)</code>	link function
<code>e(m)</code>	number of binomial trials
<code>e(wtype)</code>	weight type
<code>e(wexp)</code>	weight expression
<code>e(title_fl)</code>	family–link title
<code>e(clustvar)</code>	name of cluster variable
<code>e(offset)</code>	linear offset variable
<code>e(cons)</code>	<code>noconstant</code> or not set
<code>e(hac_kernel)</code>	HAC kernel
<code>e(hac_lag)</code>	HAC lag
<code>e(vce)</code>	<code>vcetype</code> specified in <code>vce()</code>
<code>e(vcetype)</code>	title used to label Std. Err.
<code>e(opt)</code>	type of optimization
<code>e(opt1)</code>	optimization title, line 1
<code>e(opt2)</code>	optimization title, line 2
<code>e(properties)</code>	<code>b V</code>
<code>e(predict)</code>	program used to implement <code>predict</code>
<code>e(marginsok)</code>	predictions allowed by <code>margins</code>
<code>e(marginsnotok)</code>	predictions disallowed 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

e(b)	coefficient vector
e(Cns)	constraints matrix
e(V)	variance–covariance matrix of the estimators
e(V_modelbased)	model-based variance

Functions

e(sample)	marks estimation sample
-----------	-------------------------

binreg, ml 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(df)	residual degrees of freedom
e(phi)	model scale parameter
e(aic)	model AIC, if ml
e(bic)	model BIC
e(ll)	log likelihood, if ml
e(N_clust)	number of clusters
e(chi2)	χ^2
e(p)	significance of model test
e(deviance)	deviance
e(deviance_s)	scaled deviance
e(deviance_p)	Pearson deviance
e(deviance_ps)	scaled Pearson deviance
e(dispers)	dispersion
e(dispers_s)	scaled dispersion
e(dispers_p)	Pearson dispersion
e(dispers_ps)	scaled Pearson dispersion
e(vf)	factor set by vfactor(), 1 if not set
e(rank)	rank of e(V)
e(ic)	number of iterations
e(rc)	return code
e(converged)	1 if converged, 0 otherwise

Macros

e(cmd)	binreg
e(cmdline)	command as typed
e(depvar)	name of dependent variable
e(eform)	eform() option implied by or, rr, hr, or rd
e(varfunc)	program to calculate variance function
e(varfunct)	variance title
e(varfuncf)	variance function
e(link)	program to calculate link function
e(linkt)	link title
e(linkf)	link function
e(m)	number of binomial trials
e(wtype)	weight type
e(wexp)	weight expression
e(title)	title in estimation output
e(title_fl)	family–link title
e(clustvar)	name of cluster variable
e(offset)	linear offset variable
e(cons)	noconstant or not set
e(hac_kernel)	HAC kernel
e(hac_lag)	HAC lag
e(chi2type)	Wald; type of model χ^2 test
e(vce)	vctype specified in vce()

<code>e(vctype)</code>	title used to label Std. Err.
<code>e(opt)</code>	type of optimization
<code>e(opt1)</code>	optimization title, line 1
<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>	b V
<code>e(predict)</code>	program used to implement predict
<code>e(marginsok)</code>	predictions allowed by margins
<code>e(marginsnotok)</code>	predictions disallowed by margins
<code>e(asbalanced)</code>	factor variables fvset as asbalanced
<code>e(asobserved)</code>	factor variables fvset as asobserved
Matrices	
<code>e(b)</code>	coefficient vector
<code>e(Cns)</code>	constraints matrix
<code>e(iilog)</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

Let π_i be the probability of success for the i th observation, $i = 1, \dots, N$, and let $X\beta$ be the linear predictor. The link function relates the covariates of each observation to its respective probability through the linear predictor.

In logistic regression, the logit link is used:

$$\ln\left(\frac{\pi}{1 - \pi}\right) = X\beta$$

The regression coefficient β_k represents the change in the logarithm of the odds associated with a one-unit change in the value of the X_k covariate; thus $\exp(\beta_k)$ is the ratio of the odds associated with a change of one unit in X_k .

For risk differences, the identity link $\pi = X\beta$ is used. The regression coefficient β_k represents the risk difference associated with a change of one unit in X_k . When using the identity link, you can obtain fitted probabilities outside the interval (0, 1). As suggested by Wacholder, at each iteration, fitted probabilities are checked for range conditions (and put back in range if necessary). For example, if the identity link results in a fitted probability that is smaller than $1e-4$, the probability is replaced with $1e-4$ before the link function is calculated.

A similar adjustment is made for the logarithmic link, which is used for estimating the risk ratio, $\ln(\pi) = X\beta$, where $\exp(\beta_k)$ is the risk ratio associated with a change of one unit in X_k , and for the log-complement link used to estimate the probability of no disease or health, where $\exp(\beta_k)$ represents the “health ratio” associated with a change of one unit in X_k .

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](#).

References

- Cummings, P. 2009. Methods for estimating adjusted risk ratios. *Stata Journal* 9: 175–196.
- Hardin, J. W., and M. A. Cleves. 1999. sbe29: Generalized linear models: Extensions to the binomial family. *Stata Technical Bulletin* 50: 21–25. Reprinted in *Stata Technical Bulletin Reprints*, vol. 9, pp. 140–146. College Station, TX: Stata Press.
- Kleinbaum, D. G., and M. Klein. 2010. *Logistic Regression: A Self-Learning Text*. 3rd ed. New York: Springer.
- McCullagh, P., and J. A. Nelder. 1989. *Generalized Linear Models*. 2nd ed. London: Chapman & Hall/CRC.
- Wacholder, S. 1986. Binomial regression in GLIM: Estimating risk ratios and risk differences. *American Journal of Epidemiology* 123: 174–184.
- Wright, J. T., I. G. Barrison, I. G. Lewis, K. D. MacRae, E. J. Waterson, P. J. Toplis, M. G. Gordon, N. F. Morris, and I. M. Murray-Lyon. 1983. Alcohol consumption, pregnancy and low birthweight. *Lancet* 1: 663–665.

Also see

- [R] [binreg postestimation](#) — Postestimation tools for binreg
- [R] [glm](#) — Generalized linear models
- [ME] [mecloglog](#) — Multilevel mixed-effects complementary log-log regression
- [ME] [meglm](#) — Multilevel mixed-effects generalized linear model
- [ME] [melogit](#) — Multilevel mixed-effects logistic regression
- [ME] [meprobit](#) — Multilevel mixed-effects probit regression
- [MI] [estimation](#) — Estimation commands for use with mi estimate
- [U] [20 Estimation and postestimation commands](#)