

ivfprobit — Fractional probit model with continuous endogenous covariates

[Description](#)
[Options](#)
[References](#)

[Quick start](#)
[Remarks and examples](#)
[Also see](#)

[Menu](#)
[Stored results](#)

[Syntax](#)
[Methods and formulas](#)

Description

`ivfprobit` fits a model for a fractional dependent variable, such as a rate or proportion, where one or more of the covariates are endogenous. The dependent variable must be greater than or equal to 0 and less than or equal to 1. `ivfprobit` assumes all endogenous covariates are continuous.

Quick start

Fractional probit regression of `y1` on `x` and endogenous regressor `y2` that is instrumented using `z`

```
ivfprobit y1 x (y2 = z)
```

Same as above, but with endogenous regressors `y2` and `y3`, using `z1`, `z2`, and `z3` as instruments

```
ivfprobit y1 x (y2 y3 = z1 z2 z3)
```

Menu

Statistics > Endogenous covariates > Fractional probit with endogenous covariates

Syntax

```
ivfprobit depvar [varlist1] (varlist2 = varlistiv) [if] [in] [weight] [, options]
```

*varlist*₁ is the list of exogenous variables.

*varlist*₂ is the list of endogenous variables.

*varlist*_{iv} is the list of exogenous variables used with *varlist*₁ as instruments for *varlist*₂.

<i>options</i>	Description
Model	
<code>noconstant</code>	suppress constant term
<code>constraints(<i>numlist</i>)</code>	apply specified linear constraints
SE/Robust	
<code>vce(<i>vcetype</i>)</code>	<i>vcetype</i> may be <code>robust</code> , <code>cluster <i>clustvar</i></code> , <code>bootstrap</code> , or <code>jackknife</code>
Reporting	
<code>level(#)</code>	set confidence level; default is <code>level(95)</code>
<code>first</code>	report first-stage regression
<code>nocnsreport</code>	do not display constraints
<code>display_options</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>maximize_options</code>	control the maximization process
<code>coeflegend</code>	display legend instead of statistics

*varlist*₁ and *varlist*_{iv} may contain factor variables; see [U] 11.4.3 Factor variables.

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

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

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`, `iwweights`, 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.

Options

Model

`noconstant`, `constraints(numlist)`; see [R] [Estimation options](#).

SE/Robust

`vce(vcetype)` 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`), and that use bootstrap or jackknife methods (`bootstrap`, `jackknife`); see [R] [vce_option](#).

Reporting

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

`first` requests that the parameters for the reduced-form equations showing the relationships between the endogenous variables and instruments be displayed. The default is not to show these parameter estimates.

`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](#).

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

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

Remarks and examples

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

Remarks are presented under the following headings:

[Model setup](#)
[Model identification](#)

Model setup

`ivfprobit` fits models for fractional dependent variables when one or more of the covariates is endogenous. Fractional variables can take any value in the interval $[0, 1]$; thus, `ivfprobit` is useful for modeling outcomes such as rates and proportions.

Formally, we can write the model fit by `ivfprobit` as

$$E(y_{1i} | \mathbf{x}_{1i}, \mathbf{x}_{2i}, \mathbf{y}_{2i}, u_i) = \Phi(\mathbf{y}_{2i}\boldsymbol{\beta} + \mathbf{x}_{1i}\boldsymbol{\gamma} + u_i) \quad (1)$$

$$\mathbf{y}_{2i} = \mathbf{x}_{i1}\boldsymbol{\Pi}_1 + \mathbf{x}_{i2}\boldsymbol{\Pi}_2 + v_i$$

where subscript i denotes the observation, \mathbf{y}_{2i} is a $1 \times p$ vector of continuous endogenous variables, \mathbf{x}_{1i} is a $1 \times k_1$ vector of exogenous covariates, and \mathbf{x}_{2i} is a $1 \times k_2$ vector of excluded instruments. Endogeneity arises from the possible correlation between u_i and v_i . The coefficients in vectors β and γ are the parameters of interest. Matrices $\mathbf{\Pi}_1$ and $\mathbf{\Pi}_2$ contain the coefficients of the first stage for the reduced-form equation.

To obtain parameter estimates, `ivfprobit` maximizes the same likelihood as `ivprobit` but does not require a binary dependent variable and does not require the joint density of the errors in the model to be specified correctly; see [R] `ivprobit` for more information. `ivfprobit` fits the model via quasiliikelihood estimation rather than maximum likelihood estimation. The key insight behind quasiliikelihood estimation is that we do not need to know the true distribution of the entire model to obtain consistent parameter estimates. In fact, the only requirement is the correct specification of the conditional mean given in (1) after integrating over u_i . Specifying the full distribution of the model correctly is required only if we want to obtain asymptotically efficient standard errors from maximum likelihood estimation.

`ivfprobit` does not assume that the true model is a probit model that accounts for endogeneity, such as the one fit by `ivprobit`. Therefore, the standard errors provided by maximum likelihood estimation are not appropriate. Instead, `ivfprobit` takes the maximum quasiliikelihood approach and reports robust standard errors by default.

For further discussion on quasiliikelihood estimation in the context of fractional regression, see Papke and Wooldridge (1996) and Wooldridge (2010).

► Example 1

We use a corporate 401(k) participation dataset and fit a fractional probit regression of the 401(k) participation rate (`prate`), on corporate employment size (`ltemp`) and its square, an indicator of whether the 401(k) is the sole pension plan (`sole`), and plan matching rate (`mrate`). The plan matching rate is endogenous and is instrumented using the age of the plan (`age`) and its square.

Our outcome variable `prate` has values between 0 and 1, including 1,351 firms with participation rates of 1. We assume that the functional form of the expected participation rate, after integrating over u_i , is a cumulative normal density as in (1).

We use `ivfprobit` to fit the fractional probit model, accounting for endogeneity of `mrate`.

```
. use https://www.stata-press.com/data/r18/401k
(Firm-level data on 401k participation)
. ivfprobit prate c.ltotemp##c.ltotemp i.sole (mrate = c.age##c.age)
Fitting exogenous fractional probit model:
Iteration 0: Log pseudolikelihood = -1769.7046
Iteration 1: Log pseudolikelihood = -1675.4223
Iteration 2: Log pseudolikelihood = -1674.7663
Iteration 3: Log pseudolikelihood = -1674.7661
Iteration 4: Log pseudolikelihood = -1674.7661
Fitting full model:
Iteration 0: Log pseudolikelihood = -3712.498
Iteration 1: Log pseudolikelihood = -3712.4767
Iteration 2: Log pseudolikelihood = -3712.4767
Fractional probit model with endogenous regressors
Number of obs = 4,075
Wald chi2(4) = 907.06
Prob > chi2 = 0.0000
Log pseudolikelihood = -3712.4767
```

	Coefficient	Robust std. err.	z	P> z	[95% conf. interval]	
mrate	1.907922	.0946094	20.17	0.000	1.722491	2.093353
ltotemp	-.4229273	.0744177	-5.68	0.000	-.5687833	-.2770713
c.ltotemp#						
c.ltotemp	.0217492	.0046476	4.68	0.000	.01264	.0308583
sole						
Only plan	-.1733119	.0366136	-4.73	0.000	-.2450733	-.1015504
_cons	1.904103	.3199032	5.95	0.000	1.277104	2.531102
corr(e.mrate, e.prate)	-.5690386	.0431738			-.6476498	-.4784406
sd(e.mrate)	.3989664	.0061807			.3870345	.4112661

```
Wald test of exogeneity: chi2(1) = 102.40 Prob > chi2 = 0.0000
Endogenous: mrate
Exogenous: ltotemp c.ltotemp#c.ltotemp 1.sole age c.age#c.age
```

We find a positive effect of `mrate` on the participation rate. Additionally, we see that the correlation between the unobservables, `corr(e.mrate, e.prate)`, is different from zero. This means there is evidence to support our endogeneity conjecture.

◀

Model identification

As in the basic linear instrumental-variables model, the order condition for identification requires that the number of excluded exogenous variables (that is, the additional instruments) be at least as great as the number of included endogenous variables ($k_2 \geq p$). `ivfprobit` checks this for you and issues an error message if the order condition is not met.

`ivfprobit`, like `probit` and `ivprobit`, checks the exogenous and endogenous variables to see if any of them predict the outcome variable perfectly. It will drop any offending variables and observations and then fit the model on the remaining data. Instruments that are perfect predictors do not affect estimation, so they are not checked. See [Model identification](#) in [R] [probit](#) for more information.

Stored results

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

Scalars

<code>e(N)</code>	number of observations
<code>e(k)</code>	number of 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(N_clust)</code>	number of clusters
<code>e(endog_ct)</code>	number of endogenous covariates
<code>e(p)</code>	model Wald p -value
<code>e(p_exog)</code>	exogeneity test Wald p -value
<code>e(chi2)</code>	model Wald χ^2
<code>e(chi2_exog)</code>	Wald χ^2 test of exogeneity
<code>e(rank)</code>	rank of <code>e(V)</code>
<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>ivfprobit</code>
<code>e(cmdline)</code>	command as typed
<code>e(depvar)</code>	name of dependent variable
<code>e(endog)</code>	names of endogenous variables
<code>e(exog)</code>	names of exogenous 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(chi2type)</code>	Wald; type of model χ^2 test
<code>e(vce)</code>	<i>vce</i> type specified in <code>vce()</code>
<code>e(vcetype)</code>	title used to label Std. err.
<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</code> <code>V</code>
<code>e(estat_cmd)</code>	program used to implement <code>estat</code>
<code>e(predict)</code>	program used to implement <code>predict</code>
<code>e(footnote)</code>	program used to implement the footnote display
<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(Sigma)</code>	$\widehat{\Sigma}$
<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
------------------------	-------------------------

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, 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

See *Methods and formulas* in [R] [ivprobit](#).

References

Papke, L. E., and J. M. Wooldridge. 1996. Econometric methods for fractional response variables with an application to 401(k) plan participation rates. *Journal of Applied Econometrics* 11: 619–632. [https://doi.org/10.1002/\(SICI\)1099-1255\(199611\)11:6<619::AID-JAE418>3.0.CO;2-1](https://doi.org/10.1002/(SICI)1099-1255(199611)11:6<619::AID-JAE418>3.0.CO;2-1).

Wooldridge, J. M. 2010. *Econometric Analysis of Cross Section and Panel Data*. 2nd ed. Cambridge, MA: MIT Press.

Also see

[R] [ivfprobit postestimation](#) — Postestimation tools for `ivfprobit`

[R] [fracreg](#) — Fractional response regression

[R] [gmm](#) — Generalized method of moments estimation

[R] [ivprobit](#) — Probit model with continuous endogenous covariates

[R] [ivregress](#) — Single-equation instrumental-variables regression

[R] [ivtobit](#) — Tobit model with continuous endogenous covariates

[R] [probit](#) — Probit regression

[ERM] [eprobit](#) — Extended probit regression

[SVY] [svy estimation](#) — Estimation commands for survey data

[XT] [xtprobit](#) — Random-effects and population-averaged probit models

[U] [20 Estimation and postestimation commands](#)

Stata, Stata Press, and Mata are registered trademarks of StataCorp LLC. Stata and Stata Press are registered trademarks with the World Intellectual Property Organization of the United Nations. Other brand and product names are registered trademarks or trademarks of their respective companies. Copyright © 1985–2023 StataCorp LLC, College Station, TX, USA. All rights reserved.

