

fmm — Finite mixture models using the fmm prefix

Description	Quick start	Menu
Syntax	Options	Remarks and examples
Stored results	Methods and formulas	Also see

Description

The `fmm` prefix fits finite mixture models; see [\[FMM\] fmm estimation](#) for the list of supported commands.

Quick start

Mixture of three normal distributions of `y`

```
fmm 3: regress y
```

Mixture of three linear regression models of `y` on `x1` and `x2`

```
fmm 3: regress y x1 x2
```

As above, but with class probabilities depending on `z1` and `z2`

```
fmm 3, lcp(rob(z1 z2)): regress y x1 x2
```

As above, but with additional class-specific regression covariates `x3`, `x4`, and `x5`

```
fmm, lcp(rob(z1 z2)): (regress y x1 x2 x3)    ///
                    (regress y x1 x2 x4)    ///
                    (regress y x1 x2 x5)
```

As above, but with additional class-specific probability covariates `z3` and `z4`

```
fmm: (regress y x1 x2 x3)                    ///
      (regress y x1 x2 x4, lcp(rob(z1 z2 z3)))  ///
      (regress y x1 x2 x5, lcp(rob(z1 z2 z4)))
```

Menu

Statistics > FMM (finite mixture models) > General estimation and regression

Syntax

Standard syntax

```
fmm # [if] [in] [weight] [, fmmopts] : component
```

Hybrid syntax

```
fmm [if] [in] [weight] [, fmmopts] : (component1) (component2) ...
```

where the standard syntax for *component* is

```
model depvar indepvars [, options]
```

the hybrid syntax for *component* is

```
model depvar indepvars [, lcprob(varlist) options]
```

model is an estimation command, and *options* are *model*-specific estimation options.

<i>fmmopts</i>	Description
Model	
<u>l</u> cinvariant(<i>pclassname</i>)	specify parameters that are equal across classes; default is <code>lcinvariant(none)</code>
<u>l</u> cprob(<i>varlist</i>)	specify independent variables for class probabilities
<u>l</u> clabel(<i>name</i>)	name of the categorical latent variable; default is <code>lclabel(Class)</code>
<u>l</u> cbase(#)	base latent class
<u>c</u> onstraints(<i>constraints</i>)	apply specified linear constraints
<u>c</u> ollinear	keep collinear variables
SE/Robust	
<u>v</u> ce(<i>vcetype</i>)	<i>vcetype</i> may be <code>oim</code> , <u>r</u> obust, or <u>c</u> luster <i>clustvar</i>
Reporting	
<u>l</u> evel(#)	set confidence level; default is <code>level(95)</code>
<u>n</u> ocnsreport	do not display constraints
<u>n</u> oheader	do not display header above parameter table
<u>n</u> odvheader	do not display dependent variables information in the header
notable	do not display parameter table
<i>display_options</i>	control columns and column formats, row spacing, line width, display of omitted variables and base and empty cells, and factor-variable labeling
Maximization	
<i>maximize_options</i>	control the maximization process
<u>s</u> tartvalues(<i>svmethod</i>)	method for obtaining starting values; default is <code>startvalues(factor)</code>
<u>e</u> mopts(<i>maxopts</i>)	control EM algorithm for improved starting values
<u>n</u> oestimate	do not fit the model; show starting values instead
<u>c</u> oeflegend	display legend instead of statistics

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

by, *statsby*, and *svy* are allowed; see [U] 11.1.10 **Prefix commands**.

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>pclassname</i>	Description
<code>cons</code>	intercepts and cutpoints
<code>coef</code>	fixed coefficients
<code>errvar</code>	covariances of errors
<code>scale</code>	scaling parameters
<code>all</code>	all the above
<code>none</code>	none of the above; the default

Options

Model

`lcinvariant(pclassname)` specifies which parameters of the model are constrained to be equal across the latent classes; the default is `lcinvariant(none)`.

`lcprob(varlist)` specifies that the linear prediction for a given latent class probability include the variables in *varlist*. `lcinvariant()` has no effect on these parameters.

In the standard syntax, *varlist* is used in the linear prediction for each latent class probability.

In the hybrid syntax, specify `lcprob(varlisti)` in *component*_{*i*} to include *varlist*_{*i*} in the linear prediction for the *i*th latent class probability. `lcprob()` is not allowed to be specified in *fmmopts* if it is being used in one or more *component* specifications.

In the hybrid syntax, if you specify `lcprob()` in the component that corresponds with the base latent class, the option is ignored.

`lclabel(name)` specifies a name for the categorical latent variable; the default is `lclabel(Class)`.

`lcbase(#)` specifies that # is to be treated as the base latent class.

In the standard syntax, the default is `lcbase(1)`.

In the hybrid syntax, the default base is the latent class corresponding to the first *component* that does not have `lcprob()` specified. If all components have `lcprob()`, the first *component* is the base and the `lcprob()` option specified for the first *component* is ignored.

`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`), that are robust to some kinds of misspecification (`robust`), and that allow for intragroup correlation (`cluster clustvar`); see [R] **vce_option**.

Reporting

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

`nocnsreport` suppresses the display of the constraints. Fixed-to-zero constraints that are automatically set by `fmm` are not shown in the report to keep the output manageable.

`noheader` suppresses the header above the parameter table, the display that reports the final log-likelihood value, number of observations, etc.

`nodvheader` suppresses the dependent variables information from the header above each parameter table.

`notable` suppresses the parameter tables.

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`, `tolerance(#)`, `ltolerance(#)`, `nrtolerance(#)`, `nonrtolerance`, and `from(init_specs)`; see [R] [maximize](#). These options are seldom used.

`startvalues()` specifies how starting values are to be computed. Starting values specified in `from()` override the computed starting values.

`startvalues(factor [, maxopts])` specifies that starting values are computed by assigning each observation to an initial latent class that is determined by running a `factor` analysis on all the observed variables in the specified model. This is the default.

`startvalues(classid varname [, maxopts])` specifies that starting values are computed by assigning each observation to an initial latent class specified in *varname*. *varname* is required to have each class represented in the estimation sample.

`startvalues(classpr varlist [, maxopts])` specifies that starting values are computed using the initial class probabilities specified in *varlist*. *varlist* is required to contain *g* variables for a model with *g* latent classes. The values in *varlist* are normalized to sum to 1 within each observation.

`startvalues(randomid [, draws(#) seed(#) maxopts])` specifies that starting values are computed by randomly assigning observations to initial classes.

`startvalues(randompr [, draws(#) seed(#) maxopts])` specifies that starting values are computed by randomly assigning initial class probabilities.

`startvalues(jitter [#c [#v], draws(#) seed(#) maxopts])` specifies that starting values are constructed by randomly perturbing the values from a Gaussian approximation to each outcome.

#_c is the magnitude for randomly perturbing coefficients, intercepts, cutpoints, and scale parameters; the default value is 1.

#_v is the magnitude for randomly perturbing variances for Gaussian outcomes; the default value is 1.

`startvalues(zero)` specifies that starting values are to be set to 0. This option is only useful if you use `from()` to specify starting values for some parameters and want the remaining starting values to be 0.

Most starting values options have suboptions that allow for tuning the starting values calculations:

maxopts is a subset of the standard *maximize_options*: [difficult](#), [technique\(*algorithm_spec*\)](#), [iterate\(#\)](#), [\[no\]log](#), [trace](#), [gradient](#), [showstep](#), [hessian](#), [showtolerance](#), [tolerance\(#\)](#), [ltolerance\(#\)](#), and [nrtolerance\(#\)](#); see [\[R\] maximize](#).

[draws\(#\)](#) specifies the number of random draws. For [startvalues\(randomid\)](#), [startvalues\(randompr\)](#), and [startvalues\(jitter\)](#), `fmm` will generate *#* random draws and select the starting values from the draw with the best log-likelihood value from the EM iterations. The default is [draws\(1\)](#).

[seed\(#\)](#) sets the random-number seed.

[emopts\(*maxopts*\)](#) controls maximization of the log likelihood for the EM algorithm. *maxopts* is the same subset of *maximize_options* that are allowed in the [startvalues\(\)](#) option, but some of the defaults are different for the EM algorithm. The default maximum number of iterations is [iterate\(20\)](#). The default coefficient vector tolerance is [tolerance\(1e-4\)](#). The default log-likelihood tolerance is [ltolerance\(1e-6\)](#).

[noestimate](#) specifies that the model is not to be fit. Instead, starting values are to be shown (as modified by the above options if modifications were made), and they are to be shown using the [coeflegend](#) style of output. An important use of this option is before you have modified starting values at all; you can type the following:

```
. fmm ..., ... noestimate : ...
. matrix b = e(b)
. ... (modify elements of b) ...
. fmm ..., ... from(b) : ...
```

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

[coeflegend](#) displays the legend that reveals how to specify estimated coefficients in `_b[]` notation, which you are sometimes required to type when specifying postestimation commands.

Remarks and examples

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

For a general introduction to finite mixture models, see [\[FMM\] fmm intro](#). For the list of estimation commands supported by the `fmm` prefix, see [\[FMM\] fmm estimation](#).

Examples using `fmm` can be found at

- [\[FMM\] example 1a](#) — Mixture of linear regression models
- [\[FMM\] example 1b](#) — Covariates for class membership
- [\[FMM\] example 1c](#) — Testing coefficients across class models
- [\[FMM\] example 1d](#) — Component-specific covariates
- [\[FMM\] example 2](#) — Mixture of Poisson regression models
- [\[FMM\] example 3](#) — Zero-inflated models
- [\[FMM\] example 4](#) — Mixture cure models for survival data

Stored results

fmm 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_dv)</code>	number of dependent variables
<code>e(k_cat#)</code>	number of categories for the <i>#th depvar</i> , ordinal
<code>e(k_out#)</code>	number of categories for the <i>#th depvar</i> , <code>mlogit</code>
<code>e(ll)</code>	log likelihood
<code>e(N_clust)</code>	number of clusters
<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 target model converged, 0 otherwise

Macros

<code>e(cmd)</code>	<code>gsem</code>
<code>e(cmd2)</code>	<code>fmm</code>
<code>e(cmdline)</code>	command as typed
<code>e(prefix)</code>	<code>fmm</code>
<code>e(depvar)</code>	names of dependent variables
<code>e(eqnames)</code>	names of equations
<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(model#)</code>	model for the <i>#th component</i>
<code>e(offset#)</code>	offset for the <i>#th depvar</i>
<code>e(vce)</code>	<i>vctype</i> specified in <code>vce()</code>
<code>e(vctype)</code>	title used to label Std. Err.
<code>e(opt)</code>	type of optimization
<code>e(which)</code>	<code>max</code> or <code>min</code> ; whether optimizer is to perform maximization or minimization
<code>e(method)</code>	estimation method: <code>ml</code>
<code>e(ml_method)</code>	type of <code>ml</code> 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(estat_cmd)</code>	program used to implement <code>estat</code>
<code>e(predict)</code>	program used to implement <code>predict</code>
<code>e(covariates)</code>	list of covariates
<code>e(lclass)</code>	name of latent class variable
<code>e(marginsnotok)</code>	predictions not allowed by <code>margins</code>
<code>e(marginsdefault)</code>	default <code>predict()</code> specification for <code>margins</code>
<code>e(footnote)</code>	program used to implement the footnote display
<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>	parameter vector
<code>e(b_pclass)</code>	parameter class
<code>e(cat#)</code>	categories for the <i>#th depvar</i> , ordinal
<code>e(out#)</code>	outcomes for the <i>#th depvar</i> , <code>mlogit</code>
<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>	covariance matrix of the estimators
<code>e(V_modelbased)</code>	model-based variance
<code>e(lclass_k_levels)</code>	number of levels for latent class variables
<code>e(lclass_bases)</code>	base levels for latent class variables
<code>e(_N)</code>	sample size for each component

Functions

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

Methods and formulas

Methods and formulas are presented under the following headings:

The likelihood
The EM algorithm
Survey data
Predictions

The likelihood

fmm fits finite mixture models via maximum likelihood estimation. The likelihood for the specified model is derived under the assumption that, within a given latent class, each response variable is independent and identically distributed across the estimation sample. These assumptions are conditional on the latent classes and the observed exogenous variables.

The likelihood is computed by combining the conditional likelihoods from each latent class weighted by the associated latent-class probabilities. Let θ be the vector of model parameters. For a given observation, let \mathbf{y} be the vector of observed response variables, and \mathbf{x} be the vector of independent variables. Let C be the categorical latent variable with g latent classes $1, \dots, g$. The marginal likelihood for a given observation looks something like

$$\mathcal{L}_C(\theta) = \sum_{i=1}^g \pi_i f_i(\mathbf{y}|\mathbf{x}, c_i = 1, \theta)$$

where π_i is the probability for the i th latent class, $f_i(\cdot)$ is the conditional probability density function for the observed response variables in the i th latent class, and $\mathbf{c}' = (c_1, \dots, c_g)$ is the vector of latent class indicators. When $c_i = 1$, all other elements of \mathbf{c} are zero. All auxiliary parameters are fit directly without any further parameterization, so we simply acknowledge that the auxiliary parameters are among the elements of θ .

The \mathbf{y} variables are assumed to be independent, conditional on \mathbf{x} and C , so $f_i(\cdot)$ is the product of the individual conditional densities. One exception to this is when \mathbf{y} contains the outcome and endogenous covariates for `ivregress`, in which case the Gaussian responses are actually modeled using a multivariate normal density to allow for correlated errors. This one exception does not meaningfully change the following discussion, so we make no effort to represent this distinction in the formulas.

For the i th latent class with n response variables, the conditional joint density function for a given observation is

$$f_i(\mathbf{y}|\mathbf{x}, \theta) = \prod_{j=1}^n f_{ij}(y_{ij}|\mathbf{x}, \theta)$$

All estimation commands supported by fmm model the dependence of y_{ij} on \mathbf{x} through the linear prediction

$$z_{ij} = \mathbf{x}'\beta_{ij}$$

where β_{ij} is the vector of the coefficients for y_{ij} . For notational convenience, we will overload the definitions of $f_i(\cdot)$ and $f_{ij}(\cdot)$ so that they are functions of the responses and model parameters through the linear predictions $\mathbf{z}'_i = (z_{i1}, \dots, z_{in})$. Thus $f_i(\mathbf{y}|\mathbf{x}, \theta)$ is equivalently specified as $f_i(\mathbf{y}, \mathbf{z}_i, \theta)$, and $f_{ij}(y_{ij}|\mathbf{x}, \theta)$ is equivalently specified as $f_{ij}(y_{ij}, z_{ij}, \theta)$. In this new notation, the likelihood for a given observation is

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^g \pi_i \prod_{j=1}^n f_{ij}(y_{ij}, z_{ij}, \boldsymbol{\theta}) \quad (1)$$

fmm uses the multinomial logistic distribution to model the probabilities for the latent classes. For the i th latent class, the probability is given by

$$\pi_i = \Pr(c_i = 1 | \mathbf{x}) = \frac{\exp(z_i)}{\sum_{j=1}^g \exp(z_j)}$$

where the linear prediction for the i th latent class is

$$z_i = \mathbf{x}' \boldsymbol{\gamma}_i$$

and $\boldsymbol{\gamma}_i$ is the associated vector of coefficients. If the first latent class is the base level, $\boldsymbol{\gamma}_1$ is a vector of zeros so that $z_1 = 0$ and $\exp(z_1) = 1$.

The vector $\boldsymbol{\theta}$ is therefore the set of unique model parameters taken from the following:

$\boldsymbol{\gamma}_i$ is the vector of coefficients for the i th latent class.

$\boldsymbol{\beta}_{ij}$ is the vector of coefficients for y_{ij} .

Auxiliary parameters are parameters that result from some of the distribution families.

Each latent class will have its own set of these parameters.

The EM algorithm

fmm uses the EM algorithm to refine starting values before maximizing the likelihood in (1).

The EM algorithm uses the complete-data likelihood, a likelihood where it is as if we have observed values for the latent class indicator variables \mathbf{c} . In the complete-data case, the likelihood for a given observation is

$$L(\boldsymbol{\theta}) = \prod_{i=1}^g \{\pi_i f_i(\mathbf{y}, \mathbf{z}_i, \boldsymbol{\theta})\}^{c_i}$$

so the complete-data log likelihood is

$$\log L(\boldsymbol{\theta}) = \sum_{i=1}^g c_i \{ \log \pi_i + \log f_i(\mathbf{y}, \mathbf{z}_i, \boldsymbol{\theta}) \}$$

We intend to maximize the expected complete-data log likelihood given the observed variables \mathbf{y} and \mathbf{x} . This is an iterative process in which we use the k th guess of the model parameters, denoted $\boldsymbol{\theta}_{(k)}$, then compute the next guess, $\boldsymbol{\theta}_{(k+1)}$.

In the expectation (E) step, we derive the functional form of the expected complete-data log likelihood. The complete-data log likelihood is a linear function of the latent class indicator variables, so

$$E(c_i | \mathbf{y}, \mathbf{x}, \boldsymbol{\theta}_{(k)}) = \frac{\pi_i f_i(\mathbf{y}, \mathbf{z}_i, \boldsymbol{\theta}_{(k)})}{\sum_{j=1}^g \pi_j f_j(\mathbf{y}, \mathbf{z}_j, \boldsymbol{\theta}_{(k)})}$$

We denote this posterior probability by p_i , so the expected complete-data log likelihood for a given observation is given by

$$Q(\boldsymbol{\theta} | \boldsymbol{\theta}_{(k)}) = \sum_{i=1}^g p_i \{ \log \pi_i + \log f_i(\mathbf{y}, \mathbf{z}_i, \boldsymbol{\theta}) \}$$

Note that $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_{(k)})$ is a function of $\boldsymbol{\theta}_{(k)}$ solely through the posterior probabilities p_i .

Now that we have the conditional complete-data log likelihood, the maximization (M) step is to maximize $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_{(k)})$ with respect to $\boldsymbol{\theta}$ to find $\boldsymbol{\theta}_{(k+1)}$.

Survey data

fmm supports estimation with survey data. However, only the linearized variance estimator is supported. For details on VCEs with survey data, see [SVY] [variance estimation](#).

Predictions

The predicted mean for a given response within a latent class is computed in the standard way. For example, the predicted mean for `regress` is the linear prediction, the predicted mean for `glm` is computed by applying the link function to the linear prediction, and for `ologit`, the predicted mean for a given response level is the predicted probability for that level. For survival outcomes, the formulas for predicted means (expected values) are provided in the [Survival distributions](#) section in [SEM] [methods and formulas for gsem](#).

Let \widehat{z}_i be the linear prediction for the i th latent class. The predicted probability for the i th latent class is then given by

$$\widehat{\pi}_i = \frac{\exp(\widehat{z}_i)}{\sum_{j=1}^g \exp(\widehat{z}_j)}$$

The predicted posterior probability for the i th latent class is given by

$$\widetilde{\pi}_i = \frac{\widehat{\pi}_i f_i(\mathbf{y}, \widehat{\mathbf{z}}_i, \widehat{\boldsymbol{\theta}})}{\sum_{j=1}^g \widehat{\pi}_j f_j(\mathbf{y}, \widehat{\mathbf{z}}_j, \widehat{\boldsymbol{\theta}})}$$

Let $\widehat{\mu}_i$ be the predicted mean of response y in the i th latent class. The predicted overall mean of y , using the fitted latent class probabilities, is given by

$$\widehat{\mu} = \sum_{i=1}^g \widehat{\pi}_i \widehat{\mu}_i$$

The predicted overall mean of y , using the posterior latent class probabilities, is given by

$$\tilde{\mu} = \sum_{i=1}^g \tilde{\pi}_i \hat{\mu}_i$$

Also see

[FMM] **fmm intro** — Introduction to finite mixture models

[FMM] **fmm estimation** — Fitting finite mixture models

[FMM] **fmm postestimation** — Postestimation tools for fmm

[FMM] **Glossary**

[SVY] **svy estimation** — Estimation commands for survey data