

# Package ‘projpred’

November 9, 2022

**Encoding** UTF-8

**Title** Projection Predictive Feature Selection

**Version** 2.2.2

**Date** 2022-11-09

**Description** Performs projection predictive feature selection for generalized linear and additive models as well as for generalized linear and additive multilevel models (see Piironen, Paasiniemi and Vehtari, 2020, <[doi:10.1214/20-EJS1711](https://doi.org/10.1214/20-EJS1711)>; Catalina, Bürkner and Vehtari, 2022, <<https://proceedings.mlr.press/v151/catalina22a.html>>). The package is compatible with the 'rstanarm' and 'brms' packages, but other reference models can also be used. See the documentation as well as the package vignette for more information and examples.

**License** GPL-3 | file LICENSE

**URL** <https://mc-stan.org/projpred/>, <https://discourse.mc-stan.org>

**BugReports** <https://github.com/stan-dev/projpred/issues/>

**Depends** R (>= 3.5.0)

**Imports** methods, dplyr, loo (>= 2.0.0), rstantools (>= 2.0.0), lme4, mvtnorm, ggplot2, Rcpp, utils, magrittr, mgcv, gamm4, MASS, rlang

**Suggests** rstanarm, brms, testthat, knitr, rmarkdown, glmnet, cmdstanr, bayesplot (>= 1.5.0), nlme, optimx, posterior, parallel, foreach, iterators, doParallel, future, future.callr, doFuture

**LinkingTo** Rcpp, RcppArmadillo

**Additional\_repositories** <https://mc-stan.org/r-packages/>

**LazyData** TRUE

**RoxygenNote** 7.2.1

**VignetteBuilder** knitr, rmarkdown

**NeedsCompilation** yes

**Author** Juho Piironen [aut],  
 Markus Paasiniemi [aut],  
 Alejandro Catalina [aut],  
 Frank Weber [cre, aut],  
 Aki Vehtari [aut],  
 Jonah Gabry [ctb],  
 Marco Colombo [ctb],  
 Paul-Christian Bürkner [ctb],  
 Hamada S. Badr [ctb],  
 Brian Sullivan [ctb],  
 Sölvi Rögnvaldsson [ctb]

**Maintainer** Frank Weber <fweber144@protonmail.com>

**Repository** CRAN

**Date/Publication** 2022-11-09 11:10:02 UTC

## R topics documented:

projpred-package . . . . .	3
as.matrix.projection . . . . .	5
break_up_matrix_term . . . . .	6
cl_agg . . . . .	7
cv-indices . . . . .	8
cv_varsel . . . . .	9
df_binom . . . . .	14
df_gaussian . . . . .	14
extend_family . . . . .	15
extra-families . . . . .	15
mesquite . . . . .	16
plot.vsel . . . . .	16
pred-projection . . . . .	18
predict.refmodel . . . . .	21
print.vsel . . . . .	22
print.vselsummary . . . . .	23
project . . . . .	23
refmodel-init-get . . . . .	26
solution_terms . . . . .	32
suggest_size . . . . .	33
summary.vsel . . . . .	35
varsel . . . . .	37

**Index**

**42**

## Description

The R package **projpred** performs the projection predictive variable (or "feature") selection for various regression models. We recommend to read the README file (available with enhanced formatting [online](#)) and the main vignette (topic = "projpred", but also available [online](#)) before continuing here.

Throughout the whole package documentation, we use the term "submodel" for all kinds of candidate models onto which the reference model is projected. For custom reference models, the candidate models don't need to be actual *submodels* of the reference model, but in any case (even for custom reference models), the candidate models are always actual *submodels* of the full [formula](#) used by the search procedure. In this regard, it is correct to speak of *submodels*, even in case of a custom reference model.

The following model type abbreviations will be used at multiple places throughout the documentation: GLM (generalized linear model), GLMM (generalized linear multilevel—or "mixed"—model), GAM (generalized additive model), and GAMM (generalized additive multilevel—or "mixed"—model). Note that the term "generalized" includes the Gaussian family as well.

For the projection of the reference model onto a submodel, **projpred** currently relies on the following functions (in other words, these are the workhorse functions used by the default divergence minimizer):

- Submodel without multilevel or additive terms: An internal C++ function which basically serves the same purpose as `lm()` for the `gaussian()` family and `glm()` for all other families.
- Submodel with multilevel but no additive terms: `lme4::lmer()` for the `gaussian()` family, `lme4::glmer()` for all other families.
- Submodel without multilevel but additive terms: `mgcv::gam()`.
- Submodel with multilevel and additive terms: `gamm4::gamm4()`.

The projection of the reference model onto a submodel can be run on multiple CPU cores in parallel (across the projected draws). This is powered by the **foreach** package. Thus, any parallel (or sequential) backend compatible with **foreach** can be used, e.g., the backends from packages **doParallel**, **doMPI**, or **doFuture**. Using the global option `projpred.prll_prj_trigger`, the number of projected draws below which no parallelization is applied (even if a parallel backend is registered) can be modified. Such a "trigger" threshold exists because of the computational overhead of a parallelization which makes parallelization only useful for a sufficiently large number of projected draws. By default, parallelization is turned off, which can also be achieved by supplying `Inf` (or `NULL`) to option `projpred.prll_prj_trigger`. Note that we cannot recommend parallelizing the projection on Windows because in our experience, the parallelization overhead is larger there, causing a parallel run to take longer than a sequential run. Also note that the parallelization works well for GLMs, but for GLMMs, GAMs, and GAMMs, the fitted model objects are quite big, which—when running in parallel—may lead to an excessive memory usage which in turn may crash the R session. Thus, we currently cannot recommend the parallelization for GLMMs, GAMs, and GAMMs.

## Functions

- `init_refmodel()`, `get_refmodel()` For setting up an object containing information about the reference model, the submodels, and how the projection should be carried out. Explicit calls to `init_refmodel()` and `get_refmodel()` are only rarely needed.
- `varsel()`, `cv_varsel()` For running the *search* part and the *evaluation* part for a projection predictive variable selection, possibly with cross-validation (CV).
- `summary.vsel()`, `print.vsel()`, `plot.vsel()`, `suggest_size.vsel()`, `solution_terms.vsel()` For post-processing the results from `varsel()` and `cv_varsel()`.
- `project()` For projecting the reference model onto submodel(s). Typically, this follows the variable selection, but it can also be applied directly (without a variable selection).
- `as.matrix.projection()` For extracting projected parameter draws.
- `proj_linpred()`, `proj_predict()` For making predictions from a submodel (after projecting the reference model onto it).

## Author(s)

**Maintainer:** Frank Weber <fweber144@protonmail.com>

Authors:

- Juho Piironen <juho.t.piironen@gmail.com>
- Markus Paasiniemi
- Alejandro Catalina <alecatfel@gmail.com>
- Aki Vehtari

Other contributors:

- Jonah Gabry [contributor]
- Marco Colombo [contributor]
- Paul-Christian Bürkner [contributor]
- Hamada S. Badr [contributor]
- Brian Sullivan [contributor]
- Sölvi Rögnvaldsson [contributor]

## See Also

Useful links:

- <https://mc-stan.org/projpred/>
- <https://discourse.mc-stan.org>
- Report bugs at <https://github.com/stan-dev/projpred/issues/>

---

as.matrix.projection *Extract projected parameter draws*

---

## Description

This is the `as.matrix()` method for projection objects (returned by `project()`, possibly as elements of a list). It extracts the projected parameter draws and returns them as a matrix.

## Usage

```
## S3 method for class 'projection'
as.matrix(x, nm_scheme = "auto", ...)
```

## Arguments

<code>x</code>	An object of class <code>projection</code> (returned by <code>project()</code> , possibly as elements of a list).
<code>nm_scheme</code>	The naming scheme for the columns of the output matrix. Either <code>"auto"</code> , <code>"rstanarm"</code> , or <code>"brms"</code> , where <code>"auto"</code> chooses <code>"rstanarm"</code> or <code>"brms"</code> based on the class of the reference model fit (and uses <code>"rstanarm"</code> if the reference model fit is of an unknown class).
<code>...</code>	Currently ignored.

## Value

An  $S_{\text{prj}} \times Q$  matrix of projected draws, with  $S_{\text{prj}}$  denoting the number of projected draws and  $Q$  the number of parameters.

## Examples

```
if (requireNamespace("rstanarm", quietly = TRUE)) {
  # Data:
  dat_gauss <- data.frame(y = df_gaussian$y, df_gaussian$x)

  # The "stanreg" fit which will be used as the reference model (with small
  # values for `chains` and `iter`, but only for technical reasons in this
  # example; this is not recommended in general):
  fit <- rstanarm::stan_glm(
    y ~ X1 + X2 + X3 + X4 + X5, family = gaussian(), data = dat_gauss,
    QR = TRUE, chains = 2, iter = 500, refresh = 0, seed = 9876
  )

  # Projection onto an arbitrary combination of predictor terms (with a small
  # value for `nclusters`, but only for the sake of speed in this example;
  # this is not recommended in general):
  prj <- project(fit, solution_terms = c("X1", "X3", "X5"), nclusters = 10,
    seed = 9182)
  prjmat <- as.matrix(prj)
```

```

### For further post-processing (e.g., via packages `bayesplot` and
### `posterior`), we will here ignore the fact that clustering was used
### (due to argument `nclusters` above). CAUTION: Ignoring the clustering
### is not recommended and only shown here for demonstrative purposes. A
### better solution for the clustering case is explained below.
# If the `bayesplot` package is installed, the output from
# as.matrix.projection() can be used there. For example:
if (requireNamespace("bayesplot", quietly = TRUE)) {
  print(bayesplot::mcmc_intervals(prjmat))
}
# If the `posterior` package is installed, the output from
# as.matrix.projection() can be used there. For example:
if (requireNamespace("posterior", quietly = TRUE)) {
  prjdrws <- posterior::as_draws_matrix(prjmat)
  print(posterior::summarize_draws(
    prjdrws,
    "median", "mad", function(x) quantile(x, probs = c(0.025, 0.975))
  ))
}
### Better solution for post-processing clustered draws (e.g., via
### `bayesplot` or `posterior`): Don't ignore the fact that clustering was
### used. Instead, resample the clusters according to their weights (e.g.,
### via posterior::resample_draws()). However, this requires access to the
### cluster weights which is not implemented in `projpred` yet. This
### example will be extended as soon as those weights are accessible.
}

```

---

break\_up\_matrix\_term *Break up matrix terms*

---

## Description

Sometimes there can be terms in a formula that refer to a matrix instead of a single predictor. This function breaks up the matrix term into individual predictors to handle separately, as that is probably the intention of the user.

## Usage

```
break_up_matrix_term(formula, data)
```

## Arguments

formula	A <a href="#">formula</a> for a valid model.
data	The original data.frame with a matrix as predictor.

## Value

A list containing the expanded [formula](#) and the expanded data.frame.

cl\_agg

*Weighted averaging within clusters of parameter draws***Description**

This function aggregates  $S$  parameter draws that have been clustered into  $S_{cl}$  clusters by averaging across the draws that belong to the same cluster. This averaging can be done in a weighted fashion.

**Usage**

```
cl_agg(
  draws,
  cl = seq_len(nrow(draws)),
  wdraws = rep(1, nrow(draws)),
  eps_wdraws = 0
)
```

**Arguments**

draws	An $S \times P$ matrix of parameter draws, with $P$ denoting the number of parameters.
cl	A numeric vector of length $S$ , giving the cluster indices for the draws. Draws that should be dropped (e.g., by thinning) need to have an NA in cl.
wdraws	A numeric vector of length $S$ , giving the weights of the draws. It doesn't matter whether these are normalized (i.e., sum to 1) or not because internally, these weights are normalized to sum to 1 within each cluster. Draws that should be dropped (e.g., by thinning) can (but must not necessarily) have an NA in wdraws.
eps_wdraws	A positive numeric value (typically small) which will be used to improve numerical stability: The weights of the draws within each cluster are multiplied by $1 - \text{eps\_wdraws}$ . The default of 0 should be fine for most cases; this argument only exists to help in those cases where numerical instabilities occur (which must be detected by the user; this function will not detect numerical instabilities itself).

**Value**

An  $S_{cl} \times P$  matrix of aggregated parameter draws.

**Examples**

```
set.seed(323)
S <- 100L
P <- 3L
draws <- matrix(rnorm(S * P), nrow = S, ncol = P)
# Clustering example:
S_cl <- 10L
cl_draws <- sample.int(S_cl, size = S, replace = TRUE)
draws_cl <- cl_agg(draws, cl = cl_draws)
# Clustering example with nonconstant `wdraws`:
```

```
w_draws <- rgamma(S, shape = 4)
draws_cl <- cl_agg(draws, cl = cl_draws, wdraws = w_draws)
# Thinning example (implying constant `wdraws`):
S_th <- 50L
idxs_thin <- round(seq(1, S, length.out = S_th))
th_draws <- rep(NA, S)
th_draws[idxs_thin] <- seq_len(S_th)
draws_th <- cl_agg(draws, cl = th_draws)
```

---

cv-indices

*Create cross-validation folds*


---

### Description

These are helper functions to create cross-validation (CV) folds, i.e., to split up the indices from 1 to  $n$  into  $K$  subsets ("folds") for  $K$ -fold CV. These functions are potentially useful when creating the `cvfits` and `cvfun` arguments for `init_refmodel()`. The return value is different for these two methods, see below for details.

### Usage

```
cvfolds(n, K, seed = sample.int(.Machine$integer.max, 1))
```

```
cv_ids(
  n,
  K,
  out = c("foldwise", "indices"),
  seed = sample.int(.Machine$integer.max, 1)
)
```

### Arguments

<code>n</code>	Number of observations.
<code>K</code>	Number of folds. Must be at least 2 and not exceed $n$ .
<code>seed</code>	Pseudorandom number generation (PRNG) seed by which the same results can be obtained again if needed. Passed to argument <code>seed</code> of <code>set.seed()</code> , but can also be <code>NA</code> to not call <code>set.seed()</code> at all.
<code>out</code>	Format of the output, either "foldwise" or "indices". See below for details.

### Value

`cvfolds()` returns a vector of length  $n$  such that each element is an integer between 1 and  $K$  denoting which fold the corresponding data point belongs to. The return value of `cv_ids()` depends on the `out` argument. If `out = "foldwise"`, the return value is a list with  $K$  elements, each being a list with elements `tr` and `ts` giving the training and test indices, respectively, for the corresponding fold. If `out = "indices"`, the return value is a list with elements `tr` and `ts` each being a list with  $K$  elements giving the training and test indices, respectively, for each fold.



**Examples**

```
n <- 100
set.seed(1234)
y <- rnorm(n)
cv <- cv_ids(n, K = 5, seed = 9876)
# Mean within the test set of each fold:
cvmeans <- sapply(cv, function(fold) mean(y[fold$ts]))
```

cv\_varsel

*Variable selection with cross-validation***Description**

Run the *search* part and the *evaluation* part for a projection predictive variable selection. The search part determines the solution path, i.e., the best submodel for each submodel size (number of predictor terms). The evaluation part determines the predictive performance of the submodels along the solution path. In contrast to `varsel()`, `cv_varsel()` performs a cross-validation (CV) by running the search part with the training data of each CV fold separately (an exception is explained in section "Note" below) and running the evaluation part on the corresponding test set of each CV fold.

**Usage**

```
cv_varsel(object, ...)

## Default S3 method:
cv_varsel(object, ...)

## S3 method for class 'refmodel'
cv_varsel(
  object,
  method = NULL,
  cv_method = if (!inherits(object, "datafit")) "L00" else "kfold",
  ndraws = NULL,
  nclusters = 20,
  ndraws_pred = 400,
  nclusters_pred = NULL,
  refit_prj = !inherits(object, "datafit"),
  nterms_max = NULL,
  penalty = NULL,
  verbose = TRUE,
  nloo = NULL,
  K = if (!inherits(object, "datafit")) 5 else 10,
  lambda_min_ratio = 1e-05,
  nlambda = 150,
  thresh = 1e-06,
```

```

    regul = 1e-04,
    validate_search = TRUE,
    seed = sample.int(.Machine$integer.max, 1),
    search_terms = NULL,
    ...
)

```

## Arguments

object	An object of class <code>refmodel</code> (returned by <code>get_refmodel()</code> or <code>init_refmodel()</code> ) or an object that can be passed to argument <code>object</code> of <code>get_refmodel()</code> .
...	Arguments passed to <code>get_refmodel()</code> as well as to the divergence minimizer (during a forward search and also during the evaluation part, but the latter only if <code>refit_prj</code> is <code>TRUE</code> ).
method	The method for the search part. Possible options are "L1" for L1 search and "forward" for forward search. If <code>NULL</code> , then internally, "L1" is used, except if the reference model has multilevel or additive terms or if <code>!is.null(search_terms)</code> . See also section "Details" below.
cv_method	The CV method, either "LOO" or "kfold". In the "LOO" case, a Pareto-smoothed importance sampling leave-one-out CV (PSIS-LOO CV) is performed, which avoids refitting the reference model <code>nloo</code> times (in contrast to a standard LOO CV). In the "kfold" case, a $K$ -fold CV is performed.
ndraws	Number of posterior draws used in the search part. Ignored if <code>nclusters</code> is not <code>NULL</code> or in case of L1 search (because L1 search always uses a single cluster). If both ( <code>nclusters</code> and <code>ndraws</code> ) are <code>NULL</code> , the number of posterior draws from the reference model is used for <code>ndraws</code> . See also section "Details" below.
nclusters	Number of clusters of posterior draws used in the search part. Ignored in case of L1 search (because L1 search always uses a single cluster). For the meaning of <code>NULL</code> , see argument <code>ndraws</code> . See also section "Details" below.
ndraws_pred	Only relevant if <code>refit_prj</code> is <code>TRUE</code> . Number of posterior draws used in the evaluation part. Ignored if <code>nclusters_pred</code> is not <code>NULL</code> . If both ( <code>nclusters_pred</code> and <code>ndraws_pred</code> ) are <code>NULL</code> , the number of posterior draws from the reference model is used for <code>ndraws_pred</code> . See also section "Details" below.
nclusters_pred	Only relevant if <code>refit_prj</code> is <code>TRUE</code> . Number of clusters of posterior draws used in the evaluation part. For the meaning of <code>NULL</code> , see argument <code>ndraws_pred</code> . See also section "Details" below.
refit_prj	A single logical value indicating whether to fit the submodels along the solution path again ( <code>TRUE</code> ) or to retrieve their fits from the search part ( <code>FALSE</code> ) before using those (re-)fits in the evaluation part.
nterms_max	Maximum number of predictor terms until which the search is continued. If <code>NULL</code> , then $\min(19, D)$ is used where $D$ is the number of terms in the reference model (or in <code>search_terms</code> , if supplied). Note that <code>nterms_max</code> does not count the intercept, so use <code>nterms_max = 0</code> for the intercept-only model. (Correspondingly, $D$ above does not count the intercept.)
penalty	Only relevant for L1 search. A numeric vector determining the relative penalties or costs for the predictors. A value of $0$ means that those predictors have no cost

	and will therefore be selected first, whereas Inf means those predictors will never be selected. If NULL, then 1 is used for each predictor.
verbose	A single logical value indicating whether to print out additional information during the computations.
nloo	<b>Caution:</b> Still experimental. Only relevant if <code>cv_method = "L00"</code> . Number of subsampled LOO CV folds, i.e., number of observations used for the LOO CV (anything between 1 and the original number of observations). Smaller values lead to faster computation but higher uncertainty in the evaluation part. If NULL, all observations are used, but for faster experimentation, one can set this to a smaller value.
K	Only relevant if <code>cv_method = "kfold"</code> and if the reference model was created with <code>cvfits</code> being NULL (which is the case for <code>get_refmodel.stanreg()</code> and <code>brms::get_refmodel.brmsfit()</code> ). Number of folds in $K$ -fold CV.
lambda_min_ratio	Only relevant for L1 search. Ratio between the smallest and largest lambda in the L1-penalized search. This parameter essentially determines how long the search is carried out, i.e., how large submodels are explored. No need to change this unless the program gives a warning about this.
nlambda	Only relevant for L1 search. Number of values in the lambda grid for L1-penalized search. No need to change this unless the program gives a warning about this.
thresh	Only relevant for L1 search. Convergence threshold when computing the L1 path. Usually, there is no need to change this.
regul	A number giving the amount of ridge regularization when projecting onto (i.e., fitting) submodels which are GLMs. Usually there is no need for regularization, but sometimes we need to add some regularization to avoid numerical problems.
validate_search	Only relevant if <code>cv_method = "L00"</code> . A single logical value indicating whether to cross-validate also the search part, i.e., whether to run the search separately for each CV fold (TRUE) or not (FALSE). We strongly do not recommend setting this to FALSE, because this is known to bias the predictive performance estimates of the selected submodels. However, setting this to FALSE can sometimes be useful because comparing the results to the case where this argument is TRUE gives an idea of how strongly the variable selection is (over-)fitted to the data (the difference corresponds to the search degrees of freedom or the effective number of parameters introduced by the search).
seed	Pseudorandom number generation (PRNG) seed by which the same results can be obtained again if needed. Passed to argument <code>seed</code> of <code>set.seed()</code> , but can also be NA to not call <code>set.seed()</code> at all. Here, this seed is used for clustering the reference model's posterior draws (if <code>!is.null(nclusters)</code> or <code>!is.null(nclusters_pred)</code> ), for subsampling LOO CV folds (if <code>nloo</code> is smaller than the number of observations), for sampling the folds in $K$ -fold CV, and for drawing new group-level effects when predicting from a multilevel submodel (however, not yet in case of a GAMM).
search_terms	Only relevant for forward search. A custom character vector of predictor term blocks to consider for the search. Section "Details" below describes more pre-

cisely what "predictor term block" means. The intercept ("1") is always included internally via `union()`, so there's no difference between including it explicitly or omitting it. The default `search_terms` considers all the terms in the reference model's formula.

## Details

Arguments `ndraws`, `nclusters`, `nclusters_pred`, and `ndraws_pred` are automatically truncated at the number of posterior draws in the reference model (which is 1 for `datafits`). Using less draws or clusters in `ndraws`, `nclusters`, `nclusters_pred`, or `ndraws_pred` than posterior draws in the reference model may result in slightly inaccurate projection performance. Increasing these arguments affects the computation time linearly.

For argument `method`, there are some restrictions: For a reference model with multilevel or additive formula terms, only the forward search is available. Furthermore, argument `search_terms` requires a forward search to take effect.

L1 search is faster than forward search, but forward search may be more accurate. Furthermore, forward search may find a sparser model with comparable performance to that found by L1 search, but it may also start overfitting when more predictors are added.

An L1 search may select interaction terms before the corresponding main terms are selected. If this is undesired, choose the forward search instead.

The elements of the `search_terms` character vector don't need to be individual predictor terms. Instead, they can be building blocks consisting of several predictor terms connected by the `+` symbol. To understand how these building blocks work, it is important to know how **projpred**'s forward search works: It starts with an empty vector chosen which will later contain already selected predictor terms. Then, the search iterates over model sizes  $j \in \{1, \dots, J\}$ . The candidate models at model size  $j$  are constructed from those elements from `search_terms` which yield model size  $j$  when combined with the chosen predictor terms. Note that sometimes, there may be no candidate models for model size  $j$ . Also note that internally, `search_terms` is expanded to include the intercept ("1"), so the first step of the search (model size 1) always consists of the intercept-only model as the only candidate.

As a `search_terms` example, consider a reference model with formula  $y \sim x1 + x2 + x3$ . Then, to ensure that  $x1$  is always included in the candidate models, specify `search_terms = c("x1", "x1 + x2", "x1 + x3", "x1 + x2 + x3")`. This search would start with  $y \sim 1$  as the only candidate at model size 1. At model size 2,  $y \sim x1$  would be the only candidate. At model size 3,  $y \sim x1 + x2$  and  $y \sim x1 + x3$  would be the two candidates. At the last model size of 4,  $y \sim x1 + x2 + x3$  would be the only candidate. As another example, to exclude  $x1$  from the search, specify `search_terms = c("x2", "x3", "x2 + x3")`.

## Value

An object of class `vsel`. The elements of this object are not meant to be accessed directly but instead via helper functions (see the main vignette and [projpred-package](#)).

## Note

The case `cv_method == "LOO" && !validate_search` constitutes an exception where the search part is not cross-validated. In that case, the evaluation part is based on a PSIS-LOO CV also for the submodels.

For all PSIS-LOO CVs, **projpred** calls `loo::psis()` with `r_eff = NA`. This is only a problem if there was extreme autocorrelation between the MCMC iterations when the reference model was built. In those cases however, the reference model should not have been used anyway, so we don't expect **projpred**'s `r_eff = NA` to be a problem.

## References

Magnusson, Måns, Michael Andersen, Johan Jonasson, and Aki Vehtari. 2019. "Bayesian Leave-One-Out Cross-Validation for Large Data." In *Proceedings of the 36th International Conference on Machine Learning*, edited by Kamalika Chaudhuri and Ruslan Salakhutdinov, 97:4244–53. Proceedings of Machine Learning Research. PMLR. <https://proceedings.mlr.press/v97/magnusson19a.html>.

Vehtari, Aki, Andrew Gelman, and Jonah Gabry. 2017. "Practical Bayesian Model Evaluation Using Leave-One-Out Cross-Validation and WAIC." *Statistics and Computing* 27 (5): 1413–32. doi:10.1007/s1122201696964.

Vehtari, Aki, Daniel Simpson, Andrew Gelman, Yuling Yao, and Jonah Gabry. 2022. "Pareto Smoothed Importance Sampling." arXiv. doi:10.48550/arXiv.1507.02646.

## See Also

[varsel\(\)](#)

## Examples

```
# Note: The code from this example is not executed when called via example().
# To execute it, you have to copy and paste it manually to the console.
if (requireNamespace("rstanarm", quietly = TRUE)) {
  # Data:
  dat_gauss <- data.frame(y = df_gaussian$y, df_gaussian$x)

  # The "stanreg" fit which will be used as the reference model (with small
  # values for `chains` and `iter`, but only for technical reasons in this
  # example; this is not recommended in general):
  fit <- rstanarm::stan_glm(
    y ~ X1 + X2 + X3 + X4 + X5, family = gaussian(), data = dat_gauss,
    QR = TRUE, chains = 2, iter = 500, refresh = 0, seed = 9876
  )

  # Variable selection with cross-validation (with small values
  # for `nterms_max`, `nclusters`, and `nclusters_pred`, but only for the
  # sake of speed in this example; this is not recommended in general):
  cvvs <- cv_varsel(fit, nterms_max = 3, nclusters = 5, nclusters_pred = 10,
    seed = 5555)

  # Now see, for example, `?print.vsel`, `?plot.vsel`, `?suggest_size.vsel`,
  # and `?solution_terms.vsel` for possible post-processing functions.
}
```

---

df_binom	<i>Binomial toy example</i>
----------	-----------------------------

---

**Description**

Binomial toy example

**Usage**

df\_binom

**Format**

A simulated classification dataset containing 100 observations.

**y** response, 0 or 1.

**x** predictors, 30 in total.

**Source**

<https://web.stanford.edu/~hastie/glmnet/glmnetData/BNExample.RData>

---

df_gaussian	<i>Gaussian toy example</i>
-------------	-----------------------------

---

**Description**

Gaussian toy example

**Usage**

df\_gaussian

**Format**

A simulated regression dataset containing 100 observations.

**y** response, real-valued.

**x** predictors, 20 in total. Mean and SD are approximately 0 and 1, respectively.

**Source**

<https://web.stanford.edu/~hastie/glmnet/glmnetData/QSEExample.RData>

---

extend_family	<i>Extend a family</i>
---------------	------------------------

---

**Description**

This function adds some internally required elements to an object of class family (see, e.g., [family\(\)](#)). It is called internally by [init\\_refmodel\(\)](#), so you will rarely need to call it yourself.

**Usage**

```
extend_family(family)
```

**Arguments**

family            An object of class family.

**Value**

The family object extended in the way needed by **projpred**.

---

extra-families	<i>Extra family objects</i>
----------------	-----------------------------

---

**Description**

Family objects not in the set of default [family](#) objects.

**Usage**

```
Student_t(link = "identity", nu = 3)
```

**Arguments**

link            Name of the link function. In contrast to the default [family](#) objects, this has to be a character string here.

nu              Degrees of freedom for the Student-*t* distribution.

**Value**

A family object analogous to those described in [family](#).

**Note**

Support for the [Student\\_t\(\)](#) family is still experimental.

---

mesquite

*Mesquite data set*

---

### Description

The mesquite bushes yields dataset from Gelman and Hill (2006) (<http://www.stat.columbia.edu/~gelman/arm/>).

### Usage

mesquite

### Format

The response variable is the total weight (in grams) of photosynthetic material as derived from actual harvesting of the bush. The predictor variables are:

**diam1** diameter of the canopy (the leafy area of the bush) in meters, measured along the longer axis of the bush.

**diam2** canopy diameter measured along the shorter axis.

**canopy height** height of the canopy.

**total height** total height of the bush.

**density** plant unit density (# of primary stems per plant unit).

**group** group of measurements (0 for the first group, 1 for the second group).

### Source

<http://www.stat.columbia.edu/~gelman/arm/examples/mesquite/mesquite.dat>

### References

Gelman, Andrew, and Jennifer Hill. 2006. *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge, UK: Cambridge University Press. doi:10.1017/CBO9780511790942.

---

plot.vsel

*Plot summary statistics of a variable selection*

---

### Description

This is the `plot()` method for `vsel` objects (returned by `varsel()` or `cv_varsel()`).



**Usage**

```
## S3 method for class 'vsel'
plot(
  x,
  nterms_max = NULL,
  stats = "elpd",
  deltas = FALSE,
  alpha = 0.32,
  baseline = if (!inherits(x$refmodel, "datafit")) "ref" else "best",
  thres_elpd = NA,
  ...
)
```

**Arguments**

x	An object of class <code>vsel</code> (returned by <code>varsel()</code> or <code>cv_varsel()</code> ).
nterms_max	Maximum submodel size for which the statistics are calculated. Using <code>NULL</code> is effectively the same as using <code>length(solution_terms(object))</code> . Note that <code>nterms_max</code> does not count the intercept, so use <code>nterms_max = 0</code> for the intercept-only model. For <code>plot.vsel()</code> , <code>nterms_max</code> must be at least 1.
stats	One or more character strings determining which performance statistics (i.e., utilities or losses) to calculate. Available statistics are: <ul style="list-style-type: none"> <li>• "elpd": (expected) sum of log predictive densities.</li> <li>• "mlpd": mean log predictive density, that is, "elpd" divided by the number of observations.</li> <li>• "mse": mean squared error.</li> <li>• "rmse": root mean squared error. For the corresponding standard error and lower and upper confidence interval bounds, bootstrapping is used.</li> <li>• "acc" (or its alias, "pctcorr"): classification accuracy (<code>binomial()</code> family only).</li> <li>• "auc": area under the ROC curve (<code>binomial()</code> family only). For the corresponding standard error and lower and upper confidence interval bounds, bootstrapping is used.</li> </ul>
deltas	If <code>TRUE</code> , the submodel statistics are estimated as differences from the baseline model (see argument <code>baseline</code> ) instead of estimating the actual values of the statistics.
alpha	A number determining the (nominal) coverage $1 - \alpha$ of the normal-approximation (or bootstrap; see argument <code>stats</code> ) confidence intervals. For example, in case of the normal approximation, <code>alpha = 0.32</code> corresponds to one-standard-error intervals (because of the coverage of 68%).
baseline	For <code>summary.vsel()</code> : Only relevant if <code>deltas</code> is <code>TRUE</code> . For <code>plot.vsel()</code> : Always relevant. Either "ref" or "best", indicating whether the baseline is the reference model or the best submodel found (in terms of <code>stats[1]</code> ), respectively.

`thres_elpd` Only relevant if `any(stats %in% c("elpd", "mlpd"))`. The threshold for the ELPD difference (taking the submodel's ELPD minus the baseline model's ELPD) above which the submodel's ELPD is considered to be close enough to the baseline model's ELPD. An equivalent rule is applied in case of the MLPD. See `suggest_size()` for a formalization. Supplying NA deactivates this.

... Arguments passed to the internal function which is used for bootstrapping (if applicable; see argument `stats`). Currently, relevant arguments are `B` (the number of bootstrap samples, defaulting to 2000) and `seed` (see `set.seed()`, defaulting to `sample.int(.Machine$integer.max, 1)`), but can also be NA to not call `set.seed()` at all).

### Details

As long as the reference model's performance is computable, it is always shown in the plot as a dashed red horizontal line. If `baseline = "best"`, the baseline model's performance is shown as a dotted black horizontal line. If `!is.na(thres_elpd)` and `any(stats %in% c("elpd", "mlpd"))`, the value supplied to `thres_elpd` (which is automatically adapted internally in case of the MLPD or `deltas = FALSE`) is shown as a dot-dashed gray horizontal line for the reference model and, if `baseline = "best"`, as a long-dashed green horizontal line for the baseline model.

### Examples

```
if (requireNamespace("rstanarm", quietly = TRUE)) {
  # Data:
  dat_gauss <- data.frame(y = df_gaussian$y, df_gaussian$x)

  # The "stanreg" fit which will be used as the reference model (with small
  # values for `chains` and `iter`, but only for technical reasons in this
  # example; this is not recommended in general):
  fit <- rstanarm::stan_glm(
    y ~ X1 + X2 + X3 + X4 + X5, family = gaussian(), data = dat_gauss,
    QR = TRUE, chains = 2, iter = 500, refresh = 0, seed = 9876
  )

  # Variable selection (here without cross-validation and with small values
  # for `nterms_max`, `nclusters`, and `nclusters_pred`, but only for the
  # sake of speed in this example; this is not recommended in general):
  vs <- varsel(fit, nterms_max = 3, nclusters = 5, nclusters_pred = 10,
              seed = 5555)
  print(plot(vs))
}
```

## Description

After the projection of the reference model onto a submodel, the linear predictors (for the original or a new dataset) based on that submodel can be calculated by `proj_linpred()`. The linear predictors can also be transformed to response scale. Furthermore, `proj_linpred()` returns the corresponding log predictive density values if the (original or new) dataset contains response values. The `proj_predict()` function draws from the predictive distributions (there is one such distribution for each observation from the original or new dataset) of the submodel that the reference model has been projected onto. If the projection has not been performed yet, both functions call `project()` internally to perform the projection. Both functions can also handle multiple submodels at once (for objects of class `vsel` or objects returned by a `project()` call to an object of class `vsel`; see `project()`).

## Usage

```
proj_linpred(
  object,
  newdata = NULL,
  offsetnew = NULL,
  weightsnew = NULL,
  filter_nterms = NULL,
  transform = FALSE,
  integrated = FALSE,
  .seed = sample.int(.Machine$integer.max, 1),
  ...
)
```

```
proj_predict(
  object,
  newdata = NULL,
  offsetnew = NULL,
  weightsnew = NULL,
  filter_nterms = NULL,
  nresample_clusters = 1000,
  .seed = sample.int(.Machine$integer.max, 1),
  ...
)
```

## Arguments

<code>object</code>	An object returned by <code>project()</code> or an object that can be passed to argument <code>object</code> of <code>project()</code> .
<code>newdata</code>	Passed to argument <code>newdata</code> of the reference model's <code>extract_model_data</code> function (see <code>init_refmodel()</code> ). Provides the predictor (and possibly also the response) data for the new (or old) observations. May also be <code>NULL</code> (see argument <code>extract_model_data</code> of <code>init_refmodel()</code> ). If not <code>NULL</code> , any NAs will trigger an error.
<code>offsetnew</code>	Passed to argument <code>orhs</code> of the reference model's <code>extract_model_data</code> function (see <code>init_refmodel()</code> ). Used to get the offsets for the new (or old) obser-

	variations.
weightsnew	Passed to argument wrhs of the reference model's <code>extract_model_data</code> function (see <code>init_refmodel()</code> ). Used to get the weights for the new (or old) observations.
filter_nterms	Only applies if object is an object returned by <code>project()</code> . In that case, <code>filter_nterms</code> can be used to filter object for only those elements (submodels) with a number of solution terms in <code>filter_nterms</code> . Therefore, needs to be a numeric vector or NULL. If NULL, use all submodels.
transform	For <code>proj_linpred()</code> only. A single logical value indicating whether the linear predictor should be transformed to response scale using the inverse-link function (TRUE) or not (FALSE).
integrated	For <code>proj_linpred()</code> only. A single logical value indicating whether the output should be averaged across the projected posterior draws (TRUE) or not (FALSE).
.seed	Pseudorandom number generation (PRNG) seed by which the same results can be obtained again if needed. Passed to argument seed of <code>set.seed()</code> , but can also be NA to not call <code>set.seed()</code> at all. Here, this seed is used for drawing new group-level effects in case of a multilevel submodel (however, not yet in case of a GAMM) and for drawing from the predictive distributions of the submodel(s) in case of <code>proj_predict()</code> . If a clustered projection was performed, then in <code>proj_predict()</code> , <code>.seed</code> is also used for drawing from the set of the projected clusters of posterior draws (see argument <code>nresample_clusters</code> ).
...	Arguments passed to <code>project()</code> if object is not already an object returned by <code>project()</code> .
nresample_clusters	For <code>proj_predict()</code> with clustered projection only. Number of draws to return from the predictive distributions of the submodel(s). Not to be confused with argument <code>nclusters</code> of <code>project()</code> : <code>nresample_clusters</code> gives the number of draws ( <i>with</i> replacement) from the set of clustered posterior draws after projection (with this set being determined by argument <code>nclusters</code> of <code>project()</code> ).

## Value

Let  $S_{\text{prj}}$  denote the number of (possibly clustered) projected posterior draws (short: the number of projected draws) and  $N$  the number of observations. (For `proj_linpred()` with `integrated = TRUE`, we have  $S_{\text{prj}} = 1$ .) Then, if the prediction is done for one submodel only (i.e., `length(nterms) == 1` || `!is.null(solution_terms)` in the call to `project()`):

- `proj_linpred()` returns a list with elements `pred` (predictions, i.e., the linear predictors, possibly transformed to response scale) and `lpd` (log predictive densities; only calculated if `newdata` is NULL or if `newdata` contains response values in the corresponding column). Both elements are  $S_{\text{prj}} \times N$  matrices.
- `proj_predict()` returns an  $S_{\text{prj}} \times N$  matrix of predictions where  $S_{\text{prj}}$  denotes `nresample_clusters` in case of clustered projection.

If the prediction is done for more than one submodel, the output from above is returned for each submodel, giving a named list with one element for each submodel (the names of this list being the numbers of solution terms of the submodels when counting the intercept, too).

**Examples**

```

if (requireNamespace("rstanarm", quietly = TRUE)) {
  # Data:
  dat_gauss <- data.frame(y = df_gaussian$y, df_gaussian$x)

  # The "stanreg" fit which will be used as the reference model (with small
  # values for `chains` and `iter`, but only for technical reasons in this
  # example; this is not recommended in general):
  fit <- rstanarm::stan_glm(
    y ~ X1 + X2 + X3 + X4 + X5, family = gaussian(), data = dat_gauss,
    QR = TRUE, chains = 2, iter = 500, refresh = 0, seed = 9876
  )

  # Projection onto an arbitrary combination of predictor terms (with a small
  # value for `nclusters`, but only for the sake of speed in this example;
  # this is not recommended in general):
  prj <- project(fit, solution_terms = c("X1", "X3", "X5"), nclusters = 10,
    seed = 9182)

  # Predictions (at the training points) from the submodel onto which the
  # reference model was projected:
  prjl <- proj_linpred(prj)
  prjp <- proj_predict(prj, .seed = 7364)
}

```

---

predict.refmodel

*Predictions or log predictive densities from a reference model*


---

**Description**

This is the `predict()` method for `refmodel` objects (returned by `get_refmodel()` or `init_refmodel()`). It offers three types of output which are all based on the reference model and new (or old) observations: Either the linear predictor on link scale, the linear predictor transformed to response scale, or the log predictive density.

**Usage**

```

## S3 method for class 'refmodel'
predict(
  object,
  newdata = NULL,
  ynew = NULL,
  offsetnew = NULL,
  weightsnew = NULL,
  type = "response",
  ...
)

```

**Arguments**

object	An object of class <code>refmodel</code> (returned by <code>get_refmodel()</code> or <code>init_refmodel()</code> ).
newdata	Passed to argument <code>newdata</code> of the reference model's <code>extract_model_data</code> function (see <code>init_refmodel()</code> ). Provides the predictor (and possibly also the response) data for the new (or old) observations. May also be <code>NULL</code> (see argument <code>extract_model_data</code> of <code>init_refmodel()</code> ). If not <code>NULL</code> , any NAs will trigger an error.
ynew	If not <code>NULL</code> , then this needs to be a vector of new (or old) response values. See section "Value" below.
offsetnew	Passed to argument <code>orhs</code> of the reference model's <code>extract_model_data</code> function (see <code>init_refmodel()</code> ). Used to get the offsets for the new (or old) observations.
weightsnew	Passed to argument <code>wrhs</code> of the reference model's <code>extract_model_data</code> function (see <code>init_refmodel()</code> ). Used to get the weights for the new (or old) observations.
type	Only relevant if <code>is.null(ynew)</code> . The scale on which the predictions are returned, either "link" or "response" (see <code>predict.glm()</code> but note that <code>predict.refmodel()</code> does not adhere to the typical R convention of a default prediction on link scale). For both scales, the predictions are averaged across the posterior draws.
...	Currently ignored.

**Details**

Argument `weightsnew` is only relevant if `!is.null(ynew)`.

**Value**

Either a vector of predictions (with the scale depending on argument `type`) or, if `!is.null(ynew)`, a vector of log predictive densities evaluated at `ynew`.

---

`print.vsel`

*Print results (summary) of variable selection*

---

**Description**

This is the `print()` method for `vsel` objects (returned by `varsel()` or `cv_varsel()`). It displays a summary of the results of the projection predictive variable selection by first calling `summary.vsel()` and then `print.vselsummary()`.

**Usage**

```
## S3 method for class 'vsel'
print(x, ...)
```

**Arguments**

- x                    An object of class `vsel` (returned by `vargel()` or `cv_vargel()`).
- ...                   Further arguments passed to `summary.vsel()` (apart from argument `digits` which is passed to `print.vselsummary()`).

**Value**

The output of `summary.vsel()` (invisible).

---

<code>print.vselsummary</code>	<i>Print summary of variable selection</i>
--------------------------------	--

---

**Description**

This is the `print()` method for summary objects created by `summary.vsel()`. It displays a summary of the results of the projection predictive variable selection.

**Usage**

```
## S3 method for class 'vselsummary'
print(x, digits = 1, ...)
```

**Arguments**

- x                    An object of class `vselsummary`.
- digits                Number of decimal places to be reported.
- ...                    Currently ignored.

**Value**

The output of `summary.vsel()` (invisible).

---

<code>project</code>	<i>Projection onto submodel(s)</i>
----------------------	------------------------------------

---

**Description**

Project the posterior of the reference model onto the parameter space of a single submodel consisting of a specific combination of predictor terms or (after variable selection) onto the parameter space of a single or multiple submodels of specific sizes.

**Usage**

```

project(
  object,
  nterms = NULL,
  solution_terms = NULL,
  refit_prj = TRUE,
  ndraws = 400,
  nclusters = NULL,
  seed = sample.int(.Machine$integer.max, 1),
  regul = 1e-04,
  ...
)

```

**Arguments**

object	An object which can be used as input to <code>get_refmodel()</code> (in particular, objects of class <code>refmodel</code> ).
nterms	Only relevant if object is of class <code>vsel</code> (returned by <code>varsel()</code> or <code>cv_varsel()</code> ). Ignored if <code>!is.null(solution_terms)</code> . Number of terms for the submodel (the corresponding combination of predictor terms is taken from object). If a numeric vector, then the projection is performed for each element of this vector. If <code>NULL</code> (and <code>is.null(solution_terms)</code> ), then the value suggested by the variable selection is taken (see function <code>suggest_size()</code> ). Note that <code>nterms</code> does not count the intercept, so use <code>nterms = 0</code> for the intercept-only model.
solution_terms	If not <code>NULL</code> , then this needs to be a character vector of predictor terms for the submodel onto which the projection will be performed. Argument <code>nterms</code> is ignored in that case. For an object which is not of class <code>vsel</code> , <code>solution_terms</code> must not be <code>NULL</code> .
refit_prj	A single logical value indicating whether to fit the submodels (again) ( <code>TRUE</code> ) or to retrieve the fitted submodels from object ( <code>FALSE</code> ). For an object which is not of class <code>vsel</code> , <code>refit_prj</code> must be <code>TRUE</code> . Note that currently, <code>refit_prj = FALSE</code> requires some caution, see GitHub issues #168 and #211.
ndraws	Only relevant if <code>refit_prj</code> is <code>TRUE</code> . Number of posterior draws to be projected. Ignored if <code>nclusters</code> is not <code>NULL</code> or if the reference model is of class <code>datafit</code> (in which case one cluster is used). If both ( <code>nclusters</code> and <code>ndraws</code> ) are <code>NULL</code> , the number of posterior draws from the reference model is used for <code>ndraws</code> . See also section "Details" below.
nclusters	Only relevant if <code>refit_prj</code> is <code>TRUE</code> . Number of clusters of posterior draws to be projected. Ignored if the reference model is of class <code>datafit</code> (in which case one cluster is used). For the meaning of <code>NULL</code> , see argument <code>ndraws</code> . See also section "Details" below.
seed	Pseudorandom number generation (PRNG) seed by which the same results can be obtained again if needed. Passed to argument <code>seed</code> of <code>set.seed()</code> , but can also be <code>NA</code> to not call <code>set.seed()</code> at all. Here, this seed is used for clustering the reference model's posterior draws (if <code>!is.null(nclusters)</code> ).



regul A number giving the amount of ridge regularization when projecting onto (i.e., fitting) submodels which are GLMs. Usually there is no need for regularization, but sometimes we need to add some regularization to avoid numerical problems.

... Arguments passed to `get_refmodel()` (if `get_refmodel()` is actually used; see argument object) as well as to the divergence minimizer (if `refit_prj` is TRUE).

### Details

Arguments `ndraws` and `nclusters` are automatically truncated at the number of posterior draws in the reference model (which is 1 for `datafits`). Using less draws or clusters in `ndraws` or `nclusters` than posterior draws in the reference model may result in slightly inaccurate projection performance. Increasing these arguments affects the computation time linearly.

Note that if `project()` is applied to output from `cv_varsel()`, then `refit_prj = FALSE` will take the results from the *full-data* search.

### Value

If the projection is performed onto a single submodel (i.e., `length(nterms) == 1 || !is.null(solution_terms)`), an object of class `projection` which is a list containing the following elements:

`dis` Projected draws for the dispersion parameter.

`kl` The Kullback-Leibler (KL) divergence from the reference model to the submodel. Note that in case of the Gaussian family, this is not the actual KL divergence but merely a proxy.

`weights` Weights for the projected draws.

`solution_terms` A character vector of the submodel's predictor terms.

`submdl` A list containing the submodel fits (one fit per projected draw).

`p_type` A single logical value indicating whether the reference model's posterior draws have been clustered for the projection (TRUE) or not (FALSE).

`refmodel` The reference model object.

If the projection is performed onto more than one submodel, the output from above is returned for each submodel, giving a list with one element for each submodel.

### Examples

```
if (requireNamespace("rstanarm", quietly = TRUE)) {
  # Data:
  dat_gauss <- data.frame(y = df_gaussian$y, df_gaussian$x)

  # The "stanreg" fit which will be used as the reference model (with small
  # values for `chains` and `iter`, but only for technical reasons in this
  # example; this is not recommended in general):
  fit <- rstanarm::stan_glm(
    y ~ X1 + X2 + X3 + X4 + X5, family = gaussian(), data = dat_gauss,
    QR = TRUE, chains = 2, iter = 500, refresh = 0, seed = 9876
  )
}
```

```

# Variable selection (here without cross-validation and with small values
# for `nterms_max`, `nclusters`, and `nclusters_pred`, but only for the
# sake of speed in this example; this is not recommended in general):
vs <- varsel(fit, nterms_max = 3, nclusters = 5, nclusters_pred = 10,
            seed = 5555)

# Projection onto the best submodel with 2 predictor terms (with a small
# value for `nclusters`, but only for the sake of speed in this example;
# this is not recommended in general):
prj_from_vs <- project(vs, nterms = 2, nclusters = 10, seed = 9182)

# Projection onto an arbitrary combination of predictor terms (with a small
# value for `nclusters`, but only for the sake of speed in this example;
# this is not recommended in general):
prj <- project(fit, solution_terms = c("X1", "X3", "X5"), nclusters = 10,
              seed = 9182)
}

```

---

refmodel-init-get

*Reference model and more general information*


---

## Description

Function `get_refmodel()` is a generic function whose methods usually call `init_refmodel()` which is the underlying workhorse (and may also be used directly without a call to `get_refmodel()`).

Both, `get_refmodel()` and `init_refmodel()`, create an object containing information needed for the projection predictive variable selection, namely about the reference model, the submodels, and how the projection should be carried out. For the sake of simplicity, the documentation may refer to the resulting object also as "reference model" or "reference model object", even though it also contains information about the submodels and the projection.

A "typical" reference model object is created by `get_refmodel.stanreg()` and `brms::get_refmodel.brmsfit()`, either implicitly by a call to a top-level function such as `project()`, `varsel()`, and `cv_varsel()` or explicitly by a call to `get_refmodel()`. All non-"typical" reference model objects will be called "custom" reference model objects.

Some arguments are for  $K$ -fold cross-validation ( $K$ -fold CV) only; see `cv_varsel()` for the use of  $K$ -fold CV in **projpred**.

## Usage

```

get_refmodel(object, ...)

## S3 method for class 'refmodel'
get_refmodel(object, ...)

## S3 method for class 'vsel'
get_refmodel(object, ...)

```

```

## Default S3 method:
get_refmodel(object, formula, family = NULL, ...)

## S3 method for class 'stanreg'
get_refmodel(object, ...)

init_refmodel(
  object,
  data,
  formula,
  family,
  ref_predfun = NULL,
  div_minimizer = NULL,
  proj_predfun = NULL,
  extract_model_data,
  cvfun = NULL,
  cvfits = NULL,
  dis = NULL,
  cvrefbuilder = NULL,
  ...
)

```

## Arguments

object	For <code>init_refmodel()</code> , an object that the functions from arguments <code>extract_model_data</code> and <code>ref_predfun</code> can be applied to, with a <code>NULL</code> object being treated specially (see section "Value" below). For <code>get_refmodel.default()</code> , an object of type <code>list</code> that (i) function <code>family()</code> can be applied to in order to retrieve the family (if argument <code>family</code> is <code>NULL</code> ) and (ii) has an element called <code>data</code> containing the original dataset (see argument <code>data</code> of <code>init_refmodel()</code> ), additionally to the properties required for <code>init_refmodel()</code> . For non-default methods of <code>get_refmodel()</code> , an object of the corresponding class.
...	For <code>get_refmodel.default()</code> and <code>get_refmodel.stanreg()</code> : arguments passed to <code>init_refmodel()</code> . For the <code>get_refmodel()</code> generic: arguments passed to the appropriate method. Else: ignored.
formula	The full formula to use for the search procedure. For custom reference models, this does not necessarily coincide with the reference model's formula. For general information on formulas in R, see <code>formula</code> . For multilevel formulas, see also package <b>lme4</b> (in particular, functions <code>lme4::lmer()</code> and <code>lme4::glmer()</code> ). For additive formulas, see also packages <b>mgcv</b> (in particular, function <code>mgcv::gam()</code> ) and <b>gamm4</b> (in particular, function <code>gamm4::gamm4()</code> ) as well as the notes in section "Formula terms" below.
family	An object of class <code>family</code> representing the observational model (i.e., the distributional family for the response) of the <i>submodels</i> . May be <code>NULL</code> for <code>get_refmodel.default()</code> in which case the family is retrieved from <code>object</code> . For custom reference models, <code>family</code> does not have to coincide with the family of the reference model (if the

	reference model possesses a formal family at all). In typical reference models, however, these families do coincide.
data	A data frame containing the data to use for the projection predictive variable selection. Any contrasts attributes of the dataset's columns are silently removed. For custom reference models, the columns of data do not necessarily have to coincide with those of the dataset used for fitting the reference model, but keep in mind that a row-subset of data is used for argument newdata of ref_predfun during $K$ -fold CV.
ref_predfun	Prediction function for the linear predictor of the reference model, including offsets (if existing). See also section "Arguments ref_predfun, proj_predfun, and div_minimizer" below. If object is NULL, ref_predfun is ignored and an internal default is used instead.
div_minimizer	A function for minimizing the Kullback-Leibler (KL) divergence from the reference model to a submodel (i.e., for performing the projection of the reference model onto a submodel). The output of div_minimizer is used, e.g., by proj_predfun's argument fits. See also section "Arguments ref_predfun, proj_predfun, and div_minimizer" below.
proj_predfun	Prediction function for the linear predictor of a submodel onto which the reference model is projected. See also section "Arguments ref_predfun, proj_predfun, and div_minimizer" below.
extract_model_data	A function for fetching some variables (response, observation weights, offsets) from the original dataset (supplied to argument data) or from a new dataset. See also section "Argument extract_model_data" below.
cvfun	For $K$ -fold CV only. A function that, given a fold indices vector, fits the reference model separately for each fold and returns the $K$ model fits as a list. Each of the $K$ model fits needs to be a list. If object is NULL, cvfun may be NULL for using an internal default. Only one of cvfits and cvfun needs to be provided (for $K$ -fold CV). Note that cvfits takes precedence over cvfun, i.e., if both are provided, cvfits is used.
cvfits	For $K$ -fold CV only. A list containing a sub-list called fits containing the $K$ model fits from which reference model structures are created. The cvfits list (i.e., the super-list) needs to have attributes K and folds: K has to be a single integer giving the number of folds and folds has to be an integer vector giving the fold indices (one fold index per observation). Each element of cvfits\$fits (i.e., each of the $K$ model fits) needs to be a list. Only one of cvfits and cvfun needs to be provided (for $K$ -fold CV). Note that cvfits takes precedence over cvfun, i.e., if both are provided, cvfits is used.
dis	A vector of posterior draws for the reference model's dispersion parameter or—more precisely—the posterior values for the reference model's parameter-conditional predictive variance (assuming that this variance is the same for all observations). May be NULL if the submodels have no dispersion parameter or if the submodels do have a dispersion parameter, but object is NULL (in which case $\emptyset$ is used for dis). Note that for the <code>gaussian()</code> family, dis is the standard deviation, not the variance.
cvrefbuilder	For $K$ -fold CV only. A function that, given a reference model fit for fold $k \in \{1, \dots, K\}$ (this model fit is the $k$ -th element of the return value of cvfun or

the  $k$ -th element of `cvfits$fits`, extended by elements omitted (containing the indices of the left-out observations in that fold) and `projpred_k` (containing the integer  $k$ ), returns an object of the same type as `init_refmodel()` does. Argument `cvrefbuilder` may be `NULL` for using an internal default: `get_refmodel()` if `object` is not `NULL` and a function calling `init_refmodel()` appropriately (with the assumption `dis = 0`) if `object` is `NULL`.

## Value

An object that can be passed to all the functions that take the reference model fit as the first argument, such as `varsel()`, `cv_varsel()`, `project()`, `proj_linpred()`, and `proj_predict()`. Usually, the returned object is of class `refmodel`. However, if `object` is `NULL`, the returned object is of class `datafit` as well as of class `refmodel` (with `datafit` being first). Objects of class `datafit` are handled differently at several places throughout this package.

## Formula terms

For additive models (still an experimental feature), only `mgcv::s()` and `mgcv::t2()` are currently supported as smooth terms. Furthermore, these need to be called without any arguments apart from the predictor names (symbols). For example, for smoothing the effect of a predictor  $x$ , only `s(x)` or `t2(x)` are allowed. As another example, for smoothing the joint effect of two predictors  $x$  and  $z$ , only `s(x, z)` or `t2(x, z)` are allowed (and analogously for higher-order joint effects, e.g., of three predictors).

## Arguments `ref_predfun`, `proj_predfun`, and `div_minimizer`

Arguments `ref_predfun`, `proj_predfun`, and `div_minimizer` may be `NULL` for using an internal default (see [projpred-package](#) for the functions used by the default divergence minimizer). Otherwise, let  $N$  denote the number of observations (in case of CV, these may be reduced to each fold),  $S_{\text{ref}}$  the number of posterior draws for the reference model's parameters, and  $S_{\text{prj}}$  the number of draws for the parameters of a submodel that the reference model has been projected onto (short: the number of projected draws). Then the functions supplied to these arguments need to have the following prototypes:

- `ref_predfun`: `ref_predfun(fit, newdata = NULL)` where:
  - `fit` accepts the reference model fit as given in argument `object` (but possibly re-fitted to a subset of the observations, as done in  $K$ -fold CV).
  - `newdata` accepts either `NULL` (for using the original dataset, typically stored in `fit`) or data for new observations (at least in the form of a `data.frame`).
- `proj_predfun`: `proj_predfun(fits, newdata)` where:
  - `fits` accepts a list of length  $S_{\text{prj}}$  containing this number of submodel fits. This list is the same as that returned by `project()` in its output element `submodl` (which in turn is the same as the return value of `div_minimizer`, except if `project()` was used with an object of class `vsel` based on an L1 search as well as with `refit_prj = FALSE`).
  - `newdata` accepts data for new observations (at least in the form of a `data.frame`).
- `div_minimizer` does not need to have a specific prototype, but it needs to be able to be called with the following arguments:

- `formula` accepts either a standard `formula` with a single response (if  $S_{\text{prj}} = 1$ ) or a `formula` with  $S_{\text{prj}} > 1$  response variables `cbind()`-ed on the left-hand side in which case the projection has to be performed for each of the response variables separately.
- `data` accepts a `data.frame` to be used for the projection.
- `family` accepts an object of class `family`.
- `weights` accepts either observation weights (at least in the form of a numeric vector) or `NULL` (for using a vector of ones as weights).
- `projpred_var` accepts an  $N \times S_{\text{prj}}$  matrix of predictive variances (necessary for **projpred**'s internal GLM fitter).
- `projpred_regul` accepts a single numeric value as supplied to argument `regul` of `project()`, for example.
- `...` accepts further arguments specified by the user.

The return value of these functions needs to be:

- `ref_predfun`: an  $N \times S_{\text{ref}}$  matrix.
- `proj_predfun`: an  $N \times S_{\text{prj}}$  matrix.
- `div_minimizer`: a list of length  $S_{\text{prj}}$  containing this number of submodel fits.

#### **Argument** `extract_model_data`

The function supplied to argument `extract_model_data` needs to have the prototype

```
extract_model_data(object, newdata, wrhs = NULL, orhs = NULL, extract_y = TRUE)
```

where:

- `object` accepts the reference model fit as given in argument `object` (but possibly re-fitted to a subset of the observations, as done in  $K$ -fold CV).
- `newdata` accepts either `NULL` (for using the original dataset, typically stored in `object`) or data for new observations (at least in the form of a `data.frame`).
- `wrhs` accepts at least either `NULL` (for using a vector of ones) or a right-hand side formula consisting only of the variable in `newdata` containing the weights.
- `orhs` accepts at least either `NULL` (for using a vector of zeros) or a right-hand side formula consisting only of the variable in `newdata` containing the offsets.
- `extract_y` accepts a single logical value indicating whether output element `y` (see below) shall be `NULL` (`TRUE`) or not (`FALSE`).

The return value of `extract_model_data` needs to be a list with elements `y`, `weights`, and `offset`, each being a numeric vector containing the data for the response, the observation weights, and the offsets, respectively. An exception is that `y` may also be `NULL` (depending on argument `extract_y`) or a factor.

The weights and offsets returned by `extract_model_data` will be assumed to hold for the reference model as well as for the submodels.

**Examples**

```

if (requireNamespace("rstanarm", quietly = TRUE)) {
  # Data:
  dat_gauss <- data.frame(y = df_gaussian$y, df_gaussian$x)

  # The "stanreg" fit which will be used as the reference model (with small
  # values for `chains` and `iter`, but only for technical reasons in this
  # example; this is not recommended in general):
  fit <- rstanarm::stan_glm(
    y ~ X1 + X2 + X3 + X4 + X5, family = gaussian(), data = dat_gauss,
    QR = TRUE, chains = 2, iter = 500, refresh = 0, seed = 9876
  )

  # Define the reference model explicitly:
  ref <- get_refmodel(fit)
  print(class(ref)) # gives `refmodel`
  # Now see, for example, `?varsel`, `?cv_varsel`, and `?project` for
  # possible post-processing functions. Most of the post-processing functions
  # call get_refmodel() internally at the beginning, so you will rarely need
  # to call get_refmodel() yourself.

  # A custom reference model which may be used in a variable selection where
  # the candidate predictors are not a subset of those used for the reference
  # model's predictions:
  ref_cust <- init_refmodel(
    fit,
    data = dat_gauss,
    formula = y ~ X6 + X7,
    family = gaussian(),
    extract_model_data = function(object, newdata = NULL, wrhs = NULL,
                                   orhs = NULL, extract_y = TRUE) {
      if (!extract_y) {
        resp_form <- NULL
      } else {
        resp_form <- ~ y
      }

      if (is.null(newdata)) {
        newdata <- dat_gauss
      }

      args <- projpred:::nlist(object, newdata, wrhs, orhs, resp_form)
      return(projpred:::do_call(projpred:::extract_model_data, args))
    },
    cvfun = function(folds) {
      kfold(
        fit, K = max(folds), save_fits = TRUE, folds = folds, cores = 1
      )$fits[, "fit"]
    },
    dis = as.matrix(fit)[, "sigma"]
  )
  # Now, the post-processing functions mentioned above (for example,

```

```

  # varsel(), cv_varsel(), and project()) may be applied to `ref_cust`.
}

```

---

solution_terms	<i>Retrieve predictor solution path or predictor combination</i>
----------------	--

---

### Description

This function retrieves the "solution terms" from an object. For `vsel` objects (returned by `varsel()` or `cv_varsel()`), this is the predictor solution path of the variable selection. For projection objects (returned by `project()`, possibly as elements of a list), this is the predictor combination onto which the projection was performed.

### Usage

```

solution_terms(object, ...)

## S3 method for class 'vsel'
solution_terms(object, ...)

## S3 method for class 'projection'
solution_terms(object, ...)

```

### Arguments

object	The object from which to retrieve the solution terms. Possible classes may be inferred from the names of the corresponding methods (see also the description).
...	Currently ignored.

### Value

A character vector of solution terms.

### Examples

```

if (requireNamespace("rstanarm", quietly = TRUE)) {
  # Data:
  dat_gauss <- data.frame(y = df_gaussian$y, df_gaussian$x)

  # The "stanreg" fit which will be used as the reference model (with small
  # values for `chains` and `iter`, but only for technical reasons in this
  # example; this is not recommended in general):
  fit <- rstanarm::stan_glm(
    y ~ X1 + X2 + X3 + X4 + X5, family = gaussian(), data = dat_gauss,
    QR = TRUE, chains = 2, iter = 500, refresh = 0, seed = 9876
  )
}

```



```

# Variable selection (here without cross-validation and with small values
# for `nterms_max`, `nclusters`, and `nclusters_pred`, but only for the
# sake of speed in this example; this is not recommended in general):
vs <- varsel(fit, nterms_max = 3, nclusters = 5, nclusters_pred = 10,
            seed = 5555)
print(solution_terms(vs))

# Projection onto an arbitrary combination of predictor terms (with a small
# value for `nclusters`, but only for the sake of speed in this example;
# this is not recommended in general):
prj <- project(fit, solution_terms = c("X1", "X3", "X5"), nclusters = 10,
              seed = 9182)
print(solution_terms(prj)) # gives `c("X1", "X3", "X5")`
}

```

---

suggest\_size

*Suggest submodel size*


---

## Description

This function can suggest an appropriate submodel size based on a decision rule described in section "Details" below. Note that this decision is quite heuristic and should be interpreted with caution. It is recommended to examine the results via [plot.vsel\(\)](#) and/or [summary.vsel\(\)](#) and to make the final decision based on what is most appropriate for the problem at hand.

## Usage

```

suggest_size(object, ...)

## S3 method for class 'vsel'
suggest_size(
  object,
  stat = "elpd",
  pct = 0,
  type = "upper",
  thres_elpd = NA,
  warnings = TRUE,
  ...
)

```

## Arguments

object	An object of class <code>vsel</code> (returned by <a href="#">varsel()</a> or <a href="#">cv_varsel()</a> ).
...	Arguments passed to <a href="#">summary.vsel()</a> , except for <code>object</code> , <code>stats</code> (which is set to <code>stat</code> ), <code>type</code> , and <code>deltas</code> (which is set to <code>TRUE</code> ). See section "Details" below for some important arguments which may be passed here.

stat	Performance statistic (i.e., utility or loss) used for the decision. See argument stats of <code>summary.vsel()</code> for possible choices.
pct	A number giving the proportion ( <i>not</i> percents) of the <i>relative</i> null model utility one is willing to sacrifice. See section "Details" below for more information.
type	Either "upper" or "lower" determining whether the decision is based on the upper or lower confidence interval bound, respectively. See section "Details" below for more information.
thres_elpd	Only relevant if stat %in% c("elpd", "mlpd"). The threshold for the ELPD difference (taking the submodel's ELPD minus the baseline model's ELPD) above which the submodel's ELPD is considered to be close enough to the baseline model's ELPD. An equivalent rule is applied in case of the MLPD. See section "Details" for a formalization. Supplying NA deactivates this.
warnings	Mainly for internal use. A single logical value indicating whether to throw warnings if automatic suggestion fails. Usually there is no reason to set this to FALSE.

## Details

In general (beware of special extensions below), the suggested model size is the smallest model size  $k \in \{0, 1, \dots, \text{nterms\_max}\}$  for which either the lower or upper bound (depending on argument type) of the normal-approximation (or bootstrap; see argument stat) confidence interval (with nominal coverage  $1 - \alpha$ ; see argument alpha of `summary.vsel()`) for  $U_k - U_{\text{base}}$  (with  $U_k$  denoting the  $k$ -th submodel's true utility and  $U_{\text{base}}$  denoting the baseline model's true utility) falls above (or is equal to)

$$\text{pct} \cdot (u_0 - u_{\text{base}})$$

where  $u_0$  denotes the null model's estimated utility and  $u_{\text{base}}$  the baseline model's estimated utility. The baseline model is either the reference model or the best submodel found (see argument baseline of `summary.vsel()`).

If `!is.na(thres_elpd)` and `stat = "elpd"`, the decision rule above is extended: The suggested model size is then the smallest model size  $k$  fulfilling the rule above *or*  $u_k - u_{\text{base}} > \text{thres\_elpd}$ . Correspondingly, in case of `stat = "mlpd"` (and `!is.na(thres_elpd)`), the suggested model size is the smallest model size  $k$  fulfilling the rule above *or*  $u_k - u_{\text{base}} > \frac{\text{thres\_elpd}}{N}$  with  $N$  denoting the number of observations.

For example (disregarding the special extensions in case of `stat = "elpd"` or `stat = "mlpd"`), `alpha = 0.32`, `pct = 0`, and `type = "upper"` means that we select the smallest model size for which the upper bound of the 68% confidence interval for  $U_k - U_{\text{base}}$  exceeds (or is equal to) zero, that is (if stat is a performance statistic for which the normal approximation is used, not the bootstrap), for which the submodel's utility estimate is at most one standard error smaller than the baseline model's utility estimate (with that standard error referring to the utility *difference*).

## Value

A single numeric value, giving the suggested submodel size (or NA if the suggestion failed).

The intercept is not counted by `suggest_size()`, so a suggested size of zero stands for the intercept-only model.

**Note**

Loss statistics like the root mean squared error (RMSE) and the mean squared error (MSE) are converted to utilities by multiplying them by -1, so a call such as `suggest_size(object, stat = "rmse", type = "upper")` finds the smallest model size whose upper confidence interval bound for the *negative* RMSE or MSE exceeds the cutoff (or, equivalently, has the lower confidence interval bound for the RMSE or MSE below the cutoff). This is done to make the interpretation of argument `type` the same regardless of argument `stat`.

**Examples**

```
if (requireNamespace("rstanarm", quietly = TRUE)) {
  # Data:
  dat_gauss <- data.frame(y = df_gaussian$y, df_gaussian$x)

  # The "stanreg" fit which will be used as the reference model (with small
  # values for `chains` and `iter`, but only for technical reasons in this
  # example; this is not recommended in general):
  fit <- rstanarm::stan_glm(
    y ~ X1 + X2 + X3 + X4 + X5, family = gaussian(), data = dat_gauss,
    QR = TRUE, chains = 2, iter = 500, refresh = 0, seed = 9876
  )

  # Variable selection (here without cross-validation and with small values
  # for `nterms_max`, `nclusters`, and `nclusters_pred`, but only for the
  # sake of speed in this example; this is not recommended in general):
  vs <- varsel(fit, nterms_max = 3, nclusters = 5, nclusters_pred = 10,
              seed = 5555)
  print(suggest_size(vs))
}
```

summary.vsel

*Summary statistics of a variable selection***Description**

This is the `summary()` method for `vsel` objects (returned by `varsel()` or `cv_varsel()`).

**Usage**

```
## S3 method for class 'vsel'
summary(
  object,
  nterms_max = NULL,
  stats = "elpd",
  type = c("mean", "se", "diff", "diff.se"),
  deltas = FALSE,
  alpha = 0.32,
```

```

baseline = if (!inherits(object$refmodel, "datafit")) "ref" else "best",
...
)

```

## Arguments

object	An object of class <code>vsel</code> (returned by <code>varsel()</code> or <code>cv_varsel()</code> ).
nterms_max	Maximum submodel size for which the statistics are calculated. Using NULL is effectively the same as using <code>length(solution_terms(object))</code> . Note that <code>nterms_max</code> does not count the intercept, so use <code>nterms_max = 0</code> for the intercept-only model. For <code>plot.vsel()</code> , <code>nterms_max</code> must be at least 1.
stats	One or more character strings determining which performance statistics (i.e., utilities or losses) to calculate. Available statistics are: <ul style="list-style-type: none"> <li>• "elpd": (expected) sum of log predictive densities.</li> <li>• "mlpd": mean log predictive density, that is, "elpd" divided by the number of observations.</li> <li>• "mse": mean squared error.</li> <li>• "rmse": root mean squared error. For the corresponding standard error and lower and upper confidence interval bounds, bootstrapping is used.</li> <li>• "acc" (or its alias, "pctcorr"): classification accuracy (<code>binomial()</code> family only).</li> <li>• "auc": area under the ROC curve (<code>binomial()</code> family only). For the corresponding standard error and lower and upper confidence interval bounds, bootstrapping is used.</li> </ul>
type	One or more items from "mean", "se", "lower", "upper", "diff", and "diff.se" indicating which of these to compute for each item from <code>stats</code> (mean, standard error, lower and upper confidence interval bounds, mean difference to the corresponding statistic of the reference model, and standard error of this difference, respectively). The confidence interval bounds belong to normal-approximation (or bootstrap; see argument <code>stats</code> ) confidence intervals with (nominal) coverage $1 - \alpha$ . Items "diff" and "diff.se" are only supported if <code>deltas</code> is FALSE.
deltas	If TRUE, the submodel statistics are estimated as differences from the baseline model (see argument <code>baseline</code> ) instead of estimating the actual values of the statistics.
alpha	A number determining the (nominal) coverage $1 - \alpha$ of the normal-approximation (or bootstrap; see argument <code>stats</code> ) confidence intervals. For example, in case of the normal approximation, <code>alpha = 0.32</code> corresponds to one-standard-error intervals (because of the coverage of 68%).
baseline	For <code>summary.vsel()</code> : Only relevant if <code>deltas</code> is TRUE. For <code>plot.vsel()</code> : Always relevant. Either "ref" or "best", indicating whether the baseline is the reference model or the best submodel found (in terms of <code>stats[1]</code> ), respectively.
...	Arguments passed to the internal function which is used for bootstrapping (if applicable; see argument <code>stats</code> ). Currently, relevant arguments are <code>B</code> (the number

of bootstrap samples, defaulting to 2000) and seed (see `set.seed()`, defaulting to `sample.int(.Machine$integer.max, 1)`, but can also be NA to not call `set.seed()` at all).

## Value

An object of class `vselsummary`.

## Examples

```
if (requireNamespace("rstanarm", quietly = TRUE)) {
  # Data:
  dat_gauss <- data.frame(y = df_gaussian$y, df_gaussian$x)

  # The "stanreg" fit which will be used as the reference model (with small
  # values for `chains` and `iter`, but only for technical reasons in this
  # example; this is not recommended in general):
  fit <- rstanarm::stan_glm(
    y ~ X1 + X2 + X3 + X4 + X5, family = gaussian(), data = dat_gauss,
    QR = TRUE, chains = 2, iter = 500, refresh = 0, seed = 9876
  )

  # Variable selection (here without cross-validation and with small values
  # for `nterms_max`, `nclusters`, and `nclusters_pred`, but only for the
  # sake of speed in this example; this is not recommended in general):
  vs <- varsel(fit, nterms_max = 3, nclusters = 5, nclusters_pred = 10,
              seed = 5555)
  print(summary(vs))
}
```

---

varsel

*Variable selection without cross-validation*

---

## Description

Run the *search* part and the *evaluation* part for a projection predictive variable selection. The search part determines the solution path, i.e., the best submodel for each submodel size (number of predictor terms). The evaluation part determines the predictive performance of the submodels along the solution path.

## Usage

```
varsel(object, ...)

## Default S3 method:
varsel(object, ...)

## S3 method for class 'refmodel'
```

```

varel(
  object,
  d_test = NULL,
  method = NULL,
  ndraws = NULL,
  nclusters = 20,
  ndraws_pred = 400,
  nclusters_pred = NULL,
  refit_prj = !inherits(object, "datafit"),
  nterms_max = NULL,
  verbose = TRUE,
  lambda_min_ratio = 1e-05,
  nlambdas = 150,
  thresh = 1e-06,
  regul = 1e-04,
  penalty = NULL,
  search_terms = NULL,
  seed = sample.int(.Machine$integer.max, 1),
  ...
)

```

### Arguments

object	An object of class <code>refmodel</code> (returned by <code>get_refmodel()</code> or <code>init_refmodel()</code> ) or an object that can be passed to argument <code>object</code> of <code>get_refmodel()</code> .
...	Arguments passed to <code>get_refmodel()</code> as well as to the divergence minimizer (during a forward search and also during the evaluation part, but the latter only if <code>refit_prj</code> is TRUE).
d_test	A list of the structure outlined in section "Argument d_test" below, providing test data for evaluating the predictive performance of the submodels as well as of the reference model. If NULL, the training data is used.
method	The method for the search part. Possible options are "L1" for L1 search and "forward" for forward search. If NULL, then internally, "L1" is used, except if the reference model has multilevel or additive terms or if <code>!is.null(search_terms)</code> . See also section "Details" below.
ndraws	Number of posterior draws used in the search part. Ignored if <code>nclusters</code> is not NULL or in case of L1 search (because L1 search always uses a single cluster). If both ( <code>nclusters</code> and <code>ndraws</code> ) are NULL, the number of posterior draws from the reference model is used for <code>ndraws</code> . See also section "Details" below.
nclusters	Number of clusters of posterior draws used in the search part. Ignored in case of L1 search (because L1 search always uses a single cluster). For the meaning of NULL, see argument <code>ndraws</code> . See also section "Details" below.
ndraws_pred	Only relevant if <code>refit_prj</code> is TRUE. Number of posterior draws used in the evaluation part. Ignored if <code>nclusters_pred</code> is not NULL. If both ( <code>nclusters_pred</code> and <code>ndraws_pred</code> ) are NULL, the number of posterior draws from the reference model is used for <code>ndraws_pred</code> . See also section "Details" below.

nclusters_pred	Only relevant if refit_prj is TRUE. Number of clusters of posterior draws used in the evaluation part. For the meaning of NULL, see argument ndraws_pred. See also section "Details" below.
refit_prj	A single logical value indicating whether to fit the submodels along the solution path again (TRUE) or to retrieve their fits from the search part (FALSE) before using those (re-)fits in the evaluation part.
nterms_max	Maximum number of predictor terms until which the search is continued. If NULL, then $\min(19, D)$ is used where D is the number of terms in the reference model (or in search_terms, if supplied). Note that nterms_max does not count the intercept, so use nterms_max = 0 for the intercept-only model. (Correspondingly, D above does not count the intercept.)
verbose	A single logical value indicating whether to print out additional information during the computations.
lambda_min_ratio	Only relevant for L1 search. Ratio between the smallest and largest lambda in the L1-penalized search. This parameter essentially determines how long the search is carried out, i.e., how large submodels are explored. No need to change this unless the program gives a warning about this.
nlambda	Only relevant for L1 search. Number of values in the lambda grid for L1-penalized search. No need to change this unless the program gives a warning about this.
thresh	Only relevant for L1 search. Convergence threshold when computing the L1 path. Usually, there is no need to change this.
regul	A number giving the amount of ridge regularization when projecting onto (i.e., fitting) submodels which are GLMs. Usually there is no need for regularization, but sometimes we need to add some regularization to avoid numerical problems.
penalty	Only relevant for L1 search. A numeric vector determining the relative penalties or costs for the predictors. A value of 0 means that those predictors have no cost and will therefore be selected first, whereas Inf means those predictors will never be selected. If NULL, then 1 is used for each predictor.
search_terms	Only relevant for forward search. A custom character vector of predictor term blocks to consider for the search. Section "Details" below describes more precisely what "predictor term block" means. The intercept ("1") is always included internally via union(), so there's no difference between including it explicitly or omitting it. The default search_terms considers all the terms in the reference model's formula.
seed	Pseudorandom number generation (PRNG) seed by which the same results can be obtained again if needed. Passed to argument seed of <code>set.seed()</code> , but can also be NA to not call <code>set.seed()</code> at all. Here, this seed is used for clustering the reference model's posterior draws (if !is.null(nclusters) or !is.null(nclusters_pred)) and for drawing new group-level effects when predicting from a multilevel submodel (however, not yet in case of a GAMM).

## Details

Arguments ndraws, nclusters, nclusters\_pred, and ndraws\_pred are automatically truncated at the number of posterior draws in the reference model (which is 1 for datafits). Using less

draws or clusters in `ndraws`, `nclusters`, `nclusters_pred`, or `ndraws_pred` than posterior draws in the reference model may result in slightly inaccurate projection performance. Increasing these arguments affects the computation time linearly.

For argument `method`, there are some restrictions: For a reference model with multilevel or additive formula terms, only the forward search is available. Furthermore, argument `search_terms` requires a forward search to take effect.

L1 search is faster than forward search, but forward search may be more accurate. Furthermore, forward search may find a sparser model with comparable performance to that found by L1 search, but it may also start overfitting when more predictors are added.

An L1 search may select interaction terms before the corresponding main terms are selected. If this is undesired, choose the forward search instead.

The elements of the `search_terms` character vector don't need to be individual predictor terms. Instead, they can be building blocks consisting of several predictor terms connected by the `+` symbol. To understand how these building blocks work, it is important to know how `projpred`'s forward search works: It starts with an empty vector chosen which will later contain already selected predictor terms. Then, the search iterates over model sizes  $j \in \{1, \dots, J\}$ . The candidate models at model size  $j$  are constructed from those elements from `search_terms` which yield model size  $j$  when combined with the chosen predictor terms. Note that sometimes, there may be no candidate models for model size  $j$ . Also note that internally, `search_terms` is expanded to include the intercept ("`1`"), so the first step of the search (model size 1) always consists of the intercept-only model as the only candidate.

As a `search_terms` example, consider a reference model with formula  $y \sim x_1 + x_2 + x_3$ . Then, to ensure that `x1` is always included in the candidate models, specify `search_terms = c("x1", "x1 + x2", "x1 + x3", "x1 + x2 + x3")`. This search would start with  $y \sim 1$  as the only candidate at model size 1. At model size 2,  $y \sim x_1$  would be the only candidate. At model size 3,  $y \sim x_1 + x_2$  and  $y \sim x_1 + x_3$  would be the two candidates. At the last model size of 4,  $y \sim x_1 + x_2 + x_3$  would be the only candidate. As another example, to exclude `x1` from the search, specify `search_terms = c("x2", "x3", "x2 + x3")`.

## Value

An object of class `vsel`. The elements of this object are not meant to be accessed directly but instead via helper functions (see the main vignette and [projpred-package](#)).

## Argument `d_test`

If not `NULL`, then `d_test` needs to be a `list` with the following elements:

- `data`: a `data.frame` containing the predictor variables for the test set.
- `offset`: a numeric vector containing the offset values for the test set (if there is no offset, use a vector of zeros).
- `weights`: a numeric vector containing the observation weights for the test set (if there are no observation weights, use a vector of ones).
- `y`: a numeric vector containing the response values for the test set.

## See Also

[cv\\_v`arsel`\(\)](#)



**Examples**

```
if (requireNamespace("rstanarm", quietly = TRUE)) {  
  # Data:  
  dat_gauss <- data.frame(y = df_gaussian$y, df_gaussian$x)  
  
  # The "stanreg" fit which will be used as the reference model (with small  
  # values for `chains` and `iter`, but only for technical reasons in this  
  # example; this is not recommended in general):  
  fit <- rstanarm::stan_glm(  
    y ~ X1 + X2 + X3 + X4 + X5, family = gaussian(), data = dat_gauss,  
    QR = TRUE, chains = 2, iter = 500, refresh = 0, seed = 9876  
  )  
  
  # Variable selection (here without cross-validation and with small values  
  # for `nterms_max`, `nclusters`, and `nclusters_pred`, but only for the  
  # sake of speed in this example; this is not recommended in general):  
  vs <- varsel(fit, nterms_max = 3, nclusters = 5, nclusters_pred = 10,  
              seed = 5555)  
  # Now see, for example, `?print.vsel`, `?plot.vsel`, `?suggest_size.vsel`,  
  # and `?solution_terms.vsel` for possible post-processing functions.  
}
```

# Index

**\* datasets**  
df\_binom, 14  
df\_gaussian, 14  
mesquite, 16

as.matrix(), 5  
as.matrix.projection, 5  
as.matrix.projection(), 4

binomial(), 17, 36  
break\_up\_matrix\_term, 6  
brms::get\_refmodel.brmsfit(), 11, 26

cbind(), 30  
cl\_agg, 7  
cv-indices, 8  
cv\_ids(cv-indices), 8  
cv\_ids(), 8  
cv\_varsel, 9  
cv\_varsel(), 4, 9, 16, 17, 22–26, 29, 32, 33, 35, 36, 40  
cvfolds(cv-indices), 8  
cvfolds(), 8

df\_binom, 14  
df\_gaussian, 14

extend\_family, 15  
extra-families, 15

family, 15  
family(), 15, 27  
formula, 3, 6, 27, 30

gamm4::gamm4(), 3, 27  
gaussian(), 3, 28  
get\_refmodel(refmodel-init-get), 26  
get\_refmodel(), 4, 10, 21, 22, 24–27, 29, 38  
get\_refmodel.default(), 27  
get\_refmodel.stanreg(), 11, 26, 27  
glm(), 3  
init\_refmodel(refmodel-init-get), 26  
init\_refmodel(), 4, 8, 10, 15, 19–22, 26, 27, 29, 38

lm(), 3  
lme4::glmer(), 3, 27  
lme4::lmer(), 3, 27  
loo::psis(), 13

mesquite, 16  
mgcv::gam(), 3, 27  
mgcv::s(), 29  
mgcv::t2(), 29

plot(), 16  
plot.vsel, 16  
plot.vsel(), 4, 17, 33, 36  
pred-projection, 18  
predict(), 21  
predict.glm(), 22  
predict.refmodel, 21  
predict.refmodel(), 22  
print(), 22, 23  
print.vsel, 22  
print.vsel(), 4  
print.vselsummary, 23  
print.vselsummary(), 22, 23  
proj\_linpred(pred-projection), 18  
proj\_linpred(), 4, 19, 20, 29  
proj\_predict(pred-projection), 18  
proj\_predict(), 4, 19, 20, 29  
project, 23  
project(), 4, 5, 19, 20, 25, 26, 29, 30, 32  
projpred(projpred-package), 3  
projpred-package, 3, 12, 29, 40

refmodel-init-get, 26

set.seed(), 8, 11, 18, 20, 24, 37, 39  
solution\_terms, 32  
solution\_terms.vsel(), 4

Student\_t (extra-families), 15  
Student\_t(), 15  
suggest\_size, 33  
suggest\_size(), 18, 24, 34  
suggest\_size.vsel(), 4  
summary(), 35  
summary.vsel, 35  
summary.vsel(), 4, 17, 22, 23, 33, 34, 36  
  
varsel, 37  
varsel(), 4, 9, 13, 16, 17, 22–24, 26, 29, 32,  
33, 35, 36